

Iterative Design of l_p Digital Filters

Ricardo A. Vargas and C. Sidney Burrus
Electrical and Computer Engineering Dept.
Rice University
February 2009

Abstract—The design of digital filters is a fundamental process in the context of digital signal processing. The purpose of this paper is to study the use of l_p norms (for $2 < p < \infty$) as design criteria for digital filters, and to introduce a set of algorithms for the design of *Finite (FIR)* and *Infinite (IIR) Impulse Response* digital filters based on the *Iterative Reweighted Least Squares (IRLS)* algorithm. The proposed algorithms rely on the idea of breaking the l_p filter design problem into a sequence of approximations rather than solving the original l_p problem directly. It is shown that one can efficiently design filters that arbitrarily approximate a desired l_p solution (for $2 < p < \infty$) including the commonly used l_∞ (or minimax) design problem. A method to design filters with different norms in different bands is presented (allowing the user for better control of the signal and noise behavior per band). Among the main contributions of this work is a method for the design of *magnitude* l_p IIR filters. Experimental results show that the algorithms in this work are robust and efficient, improving over traditional off-the-shelf optimization tools. The group of proposed algorithms form a flexible collection that offers robustness and efficiency for a wide variety of digital filter design applications.

I. INTRODUCTION

The design of digital filters has fundamental importance in digital signal processing. One can find applications of digital filters in many diverse areas of science and engineering including medical imaging, audio and video processing, oil exploration, and highly sophisticated military applications. Furthermore, each of these applications benefits from digital filters in particular ways, thus requiring different properties from the filters they employ. Therefore it is of critical importance to have efficient design methods that can shape filters according to the user's needs.

In this work we use the discrete l_p norm as the criterion for designing efficient digital filters. We also introduce a set of algorithms, all based on the *Iterative Reweighted Least Squares (IRLS)* method, to solve a variety of relevant digital filter design problems. The proposed family of algorithms has proven to be efficient in practice; these algorithms share theoretical justification for their use and implementation. Finally, the document makes a point about the relevance of the l_p norm as a useful tool in filter design applications.

The rest of this chapter is devoted to motivating the problem. Section I-A introduces the general filter design problem and some of the signal processing concepts relevant to this work. Section I-C presents the *basic* Iterative Reweighted Least Squares method, one of the key concepts in this document. Section I-D introduces *Finite Impulse Response (FIR)* filters and covers theoretical motivations for l_p design, including previous knowledge in l_p optimization (both from experiences in filter design as well as other fields of science and engineering). Similarly, Section I-E introduces *Infinite Impulse Response (IIR)* filters. These last two sections lay down the structure of the proposed algorithms, and provide an outline for the main contributions of this work.

Chapters II and III formally introduce the different l_p filter design problems considered in this work and discuss their IRLS-based algorithms and corresponding results. Each of these chapters provides a literary review of related previous work as well as a discussion on the proposed methods and their corresponding results.

An important contribution of this work is the extension of known and well understood concepts in l_p FIR filter design to the IIR case.

A. Digital filter design

When designing digital filters for signal processing applications one is often interested in creating objects $\mathbf{h} \in \mathbb{R}^N$ in order to alter some of the properties of a given vector $\mathbf{x} \in \mathbb{R}^M$ (where $0 < M, N < \infty$). Often the properties of \mathbf{x} that we are interested in changing lie in the frequency domain, with $\mathbf{X} = \mathcal{F}(\mathbf{x})$ being the *frequency domain* representation of \mathbf{x} given by

$$\mathbf{x} \xleftrightarrow{\mathcal{F}} \mathbf{X} = \mathbf{A}_X e^{j\omega\phi_X}$$

where \mathbf{A}_X and ϕ_X are the *amplitude* and *phase* components of \mathbf{x} , and $\mathcal{F}(\cdot) : \mathbb{R}^N \mapsto \mathbb{R}^\infty$ is the *Fourier transform* operator defined by

$$\mathcal{F}\{\mathbf{h}\} = H(\omega) \triangleq \sum_{n=0}^{N-1} h_n e^{-j\omega n} \quad \forall \omega \in [-\pi, \pi] \quad (1)$$

So the idea in filter design is to create filters \mathbf{h} such that the Fourier transform \mathbf{H} of \mathbf{h} possesses desirable *amplitude* and *phase* characteristics.

The *filtering* operator is the convolution operator $(*)$ defined by

$$(\mathbf{x} * \mathbf{h})(n) = \sum_m x(m)h(n-m)$$

An important property of the convolution operator is the *Convolution Theorem* [1] which states that

$$\mathbf{x} * \mathbf{h} \xleftrightarrow{\mathcal{F}} \mathbf{X} \cdot \mathbf{H} = (\mathbf{A}_X \cdot \mathbf{A}_H) e^{j\omega(\phi_X + \phi_H)} \quad (2)$$

where $\{\mathbf{A}_X, \phi_X\}$ and $\{\mathbf{A}_H, \phi_H\}$ represent the amplitude and phase components of \mathbf{X} and \mathbf{H} respectively. It can be seen that by *filtering* \mathbf{x} with \mathbf{h} one can apply a *scaling* operator to the amplitude of \mathbf{x} and a *biasing* operator to its phase.

A common use of digital filters is to remove a certain *band* of frequencies from the frequency spectra of \mathbf{x} (such as typical *lowpass* filters). Other types of filters include *band-pass*, *high-pass* or *band-reject* filters, depending on the range of frequencies that they alter.

B. The notion of approximation in l_p filter design

Once a filter design concept has been selected, the design problem becomes finding the optimal vector $\mathbf{h} \in \mathbb{R}^n$ that most closely *approximates* our desired frequency response concept (we will denote such *optimal* vector by \mathbf{h}^*). This approximation problem will heavily depend on the *measure* by which we evaluate all vectors $\mathbf{h} \in \mathbb{R}^N$ to choose \mathbf{h}^* .

In this document we consider the discrete l_p norms defined by

$$\|\mathbf{a}\|_p = \sqrt[p]{\sum_k |a_k|^p} \quad \forall \mathbf{a} \in \mathbb{R}^N \quad (3)$$

as measures of optimality, and consider a number of filter design problems based upon this criterion. The work explores the *Iterative*

Reweighted Least Squares (IRLS) approach as a design tool, and provides a number of algorithms based on this method. Finally, this work considers critical theoretical aspects and evaluates the numerical properties of the proposed algorithms in comparison to existing *general purpose* methods commonly used. It is the belief of the author (as well as the author's advisor) that the IRLS approach offers a more tailored route to the l_p filter design problems considered, and that it contributes an example of a *made-for-purpose* algorithm best suited to the characteristics of l_p filter design.

C. The IRLS algorithm

Iterative Reweighted Least Squares (IRLS) algorithms define a family of iterative methods that solve an otherwise complicated numerical optimization problem by breaking it into a series of *weighted least squares* (WLS) problems, each one easier in principle than the original problem. At iteration i one must solve a weighted least squares problem of the form

$$\min_{h_i} \|w(h_{i-1})f(h_i)\|_2 \quad (4)$$

where $w(\cdot)$ is a specific weighting function and $f(\cdot)$ is a function of the filter. Obviously a large class of problems could be written in this form (large in the sense that both $w(\cdot)$ and $f(\cdot)$ can be defined arbitrarily). One case worth considering is the *linear approximation* problem defined by

$$\min_h \|D - Ch\| \quad (5)$$

where $D \in \mathbb{R}^M$ and $C \in \mathbb{R}^{M \times N}$ are given, and $\|\cdot\|$ is an arbitrary measure. One could write $f(\cdot)$ in (4) as

$$f(h) = D - Ch$$

and attempt to find a suitable function $w(\cdot)$ to minimize the arbitrary norm $\|\cdot\|$ in (5). In vector notation, at iteration i one can write (4) as follows,

$$\min_{h_i} \|w(h_{i-1})(D - Ch_i)\|_2 \quad (6)$$

One can show that the solution of (6) for any iteration is given by

$$h = (C^T W C)^{-1} C^T W D$$

with $W = \text{diag}(w^2)$ (where w is the weighting vector). To solve problem (6) above, one could use the following algorithm:

- 1) Set initial weights w_0
- 2) At the i -th iteration find $h_i = (C^T W_{i-1} C)^{-1} C^T W_{i-1} D$
- 3) Update W_i as a function of h_i (i.e. $W_i = W(h_i)$)
- 4) Iterate steps 2 and 3 until a certain stopping criterion is reached

This method will be referred in this work as the *basic* IRLS algorithm.

An IRLS algorithm is said to *converge* if the algorithm produces a sequence of points h_i such that

$$\lim_{i \rightarrow \infty} h_i = h^*$$

where h^* is a *fixed point* defined by

$$h^* = (C^T W^* C)^{-1} C^T W^* D$$

with $W^* = W(h^*)$. In principle one would want $h^* = h^*$ (as defined in Section I-B).

IRLS algorithms have been used in different areas of science and engineering. Their attractiveness stem from the idea of simplifying a difficult problem as a sequence of weighted least squares problems that can be solved efficiently with programs such as Matlab or LAPACK. However (as it was mentioned above) success is determined by the existence of a weighting function that leads to a fixed point that happens to be at least a local solution of the problem in question.

This might not be the case for any given problem. In the case of l_p optimization one can justify the use of IRLS methods by means of the following theorem:

Theorem 1 (Weight Function Existence theorem): Let $g_k(\omega)$ be a Chebyshev set and define

$$H(h; \omega) = \sum_{k=0}^M h_k g_k(\omega)$$

where $h = (h_0, h_1, \dots, h_M)^T$. Then, given $D(\omega)$ continuous on $[0, \pi]$ and $1 < q < p \leq \infty$ the following are identical sets:

- $\{h \mid H(h; \omega) \text{ is a best weighted } L_p \text{ approximation to } D(\omega) \text{ on } [0, \pi]\}$.
- $\{h \mid H(h; \omega) \text{ is a best weighted } L_q \text{ approximation to } D(\omega) \text{ on } [0, \pi]\}$.

Furthermore, the theorem above is valid if the interval $[0, \pi]$ is replaced by a finite point set $\Omega \subset [0, \pi]$ (this theorem is accredited to Motzkin and Walsh [2], [3]).

Theorem 1 is fundamental since it establishes that weights exist so that the solution of an L_p problem is indeed the solution of a weighted L_q problem (for arbitrary $p, q > 1$). Furthermore the results of Theorem 1 remain valid for l_p and l_q . For our purposes, this theorem establishes the existence of a weighting function so that the solution of a weighted l_2 problem is indeed the solution of an l_p problem; the challenge then is to find the corresponding *weighting function*. The remainder of this document explores this task for a number of relevant filter design problems and provides a consistent computational framework.

D. Finite Impulse Response (FIR) l_p design

A *Finite Impulse Response (FIR)* filter is an ordered vector $h \in \mathbb{R}^N$ (where $1 \leq N < \infty$), with a complex polynomial form in the frequency domain given by

$$H(\omega) = \sum_{n=0}^{N-1} h_n e^{-j\omega n}$$

The filter $H(\omega)$ contains amplitude and phase components $\{A_H(\omega), \phi_H(\omega)\}$ that can be designed to suit the user's purpose.

Given a desired frequency response $D(\omega)$, the general l_p approximation problem is given by

$$\min_h \|D(\omega) - H(h; \omega)\|_p$$

In the most basic scenario $D(\omega)$ would be a complex valued function, and the optimization algorithm would minimize the l_p norm of the complex error function $\epsilon(\omega) = D(\omega) - H(\omega)$; we refer to this case as the *complex* l_p design problem (refer to Section II-C).

One of the caveats of solving complex approximation problems is that the user must provide desired magnitude and phase specifications. In many applications one is interested in removing or altering a range of frequencies from a signal; in such instances it might be more convenient to only provide the algorithm with a desired magnitude function while allowing the algorithm to find a phase that corresponds to the optimal magnitude design. The *magnitude* l_p design problem is given by

$$\min_h \|D(\omega) - |H(h; \omega)|\|_p$$

where $D(\omega)$ is a real, positive function. This problem is discussed in Section II-D.

Another problem that uses no phase information is the *linear phase* l_p problem. It will be shown in Section II-B that this problem can be formulated so that only real functions are involved in the optimization

problem (since the phase component of $H(\omega)$ has a specific linear form).

An interesting case results from the idea of combining different norms in different frequency bands of a desired function $D(\omega)$. One could assign different p -values for different bands (for example, minimizing the error energy (ε_2) in the passband while using a minimax error (ε_∞) approach in the stopband to keep control of noise). The *frequency-varying* l_p problem is formulated as follows,

$$\min_{\mathbf{h}} \|(D - H)(\omega_{pb})\|_p + \|(D - H)(\omega_{sb})\|_q$$

where $\{\omega_{pb}, \omega_{sb}\}$ are the passband and stopband frequency ranges respectively (and $2 < p, q < \infty$).

Perhaps the most relevant problem addressed in this work is the *Constrained Least Squares* (CLS) problem. In a continuous sense, a CLS problem is defined by

$$\begin{aligned} \min_{\mathbf{h}} \quad & \|d(\omega) - H(\omega)\|_2 \\ \text{subject to} \quad & |d(\omega) - H(\omega)| \leq \tau \end{aligned}$$

The idea is to minimize the error energy across all frequencies, but ensuring first that the error at each frequency does not exceed a given tolerance τ . Section II-F explains the details for this problem and shows that this type of formulation makes good sense in filter design and can efficiently be solved via IRLS methods.

1) *The IRLS algorithm and FIR literature review:* A common approach to dealing with highly structured approximation problems consists in breaking a complex problem into a series of simpler, smaller problems. Often, one can even prove important mathematical properties in this way. Consider the l_p approximation problem introduced in (3),

$$\min_{\mathbf{h}} \|f(\mathbf{h})\|_p \quad (7)$$

For simplicity at this point we can assume that $f(\cdot) : \mathbb{R}^N \mapsto \mathbb{R}^M$ is linear. It is relevant to mention that (7) is equivalent to

$$\min_{\mathbf{h}} \|f(\mathbf{h})\|_p^p \quad (8)$$

In its most basic form the l_p IRLS algorithm works by rewriting (8) into a weighted least squares problem of the form

$$\min_{\mathbf{h}} \|w(\mathbf{h})f(\mathbf{h})\|_2^2 \quad (9)$$

Since a linear weighted least squares problem like (9) has a closed form solution, it can be solved in one step. Then the solution is used to update the weighting function, which is kept constant for the next closed form solution and so on (as discussed in Section I-C).

One of the earlier works on the use of IRLS methods for l_p approximation was written by Charles Lawson [4]–[6], in part motivated by problems that might not have a suitable l_∞ algorithm. He looked at a basic form of the IRLS method to solve l_∞ problems and extended it by proposing a multiplicative update of the weighting coefficients at each iteration (that is, $w_{k+1}(\omega) = f(\omega) \cdot w_k(\omega)$). Lawson's method triggered a number of papers and ideas; however his method is sensitive to the weights becoming numerically zero; in this case the algorithm must restart. A number of ideas [5], [6] have been proposed (some from Lawson himself) to prevent or deal with these occurrences, and in general his method is considered somewhat slow.

John Rice and Karl Usow [5], [7] extended Lawson's method to the general l_p problem ($2 < p < \infty$) by developing an algorithm based on Lawson's that also updates the weights in a multiplicative form. They used the results from Theorem 1 by Motzkin and Walsh [2], [3] to guarantee that a solution indeed exists for the l_p problem. They defined

$$w_{k+1}(\omega) = w_k^\alpha(\omega) |\epsilon_k(\omega)|^\beta$$

where

$$\alpha = \frac{\gamma(p-2)}{\gamma(p-2)+1}$$

and

$$\beta = \frac{\alpha}{2\gamma} = \frac{p-2}{2(\gamma(p-2)+1)}$$

with γ being a convergence parameter and $\epsilon(\omega) = d(\omega) - H(\omega)$. The rest of the algorithm works the same way as the basic IRLS method; however the proper selection of γ could allow for strong convergence (note that for $\gamma = 0$ we obtain the basic IRLS algorithm).

Another approach to solve (7) consists in a *partial updating* strategy of the *filter coefficients* rather than the *weights*, by using a temporary coefficient vector defined by

$$\hat{\mathbf{a}}_{k+1} = [\mathbf{C}^T \mathbf{W}_k^T \mathbf{W}_k \mathbf{C}]^{-1} \mathbf{C}^T \mathbf{W}_k^T \mathbf{W}_k \mathbf{A}_d \quad (10)$$

The filter coefficients after each iteration are then calculated by

$$\mathbf{a}_{k+1} = \lambda \hat{\mathbf{a}}_{k+1} + (1 - \lambda) \mathbf{a}_k \quad (11)$$

where λ is a *convergence parameter* (with $0 < \lambda < 1$). This approach is known as the Karlovitz method [8], and it has been claimed that it converges to the global optimal solution for **even** values of p such that $4 \leq p < \infty$. However, in practice several convergence problems have been found even under such assumptions. One drawback is that the convergence parameter λ has to be optimized at each iteration via an expensive line search process. Therefore the overall execution time becomes rather large.

S. W. Kahng [9] developed an algorithm based on Newton-Raphson's method that uses

$$\lambda = \frac{1}{p-1} \quad (12)$$

to get

$$\mathbf{a}_{k+1} = \frac{\hat{\mathbf{a}}_{k+1} + (p-2)\mathbf{a}_k}{p-1} \quad (13)$$

This selection for λ is based upon Newton's method to minimize ϵ (the same result was derived independently by Fletcher, Grant and Hebden [10]). The rest of the algorithm follows Karlovitz's approach; however since λ is fixed there is no need to perform the linear search for its best value. Since Kahng's method is based on Newton's method, it converges quadratically to the optimal solution. Kahng proved that his method converges for all cases of λ and for any problem (at least in theory). It can be seen that Kahng's method is a particular case of Karlovitz's algorithm, with λ as defined in (12). Newton-Raphson based algorithms are not warranted to converge to the optimal solution unless they are somewhat close to the solution since they require to know and invert the Hessian matrix of the objective function (which must be *positive definite* [11]). However, their associated quadratic convergence makes them an appealing option.

Burrus, Barreto and Selesnick developed a method [7], [12], [13] that combines the powerful quadratic convergence of Newton's methods with the robust initial convergence of the basic IRLS method, thus overcoming the initial sensitivity of Newton-based algorithms and the slow linear convergence of Lawson-based methods. To accelerate initial convergence, their approach to solve (7) uses $p = \sigma * 2$, where σ is a convergence parameter (with $1 < \sigma \leq 2$). At any given iteration, p increases its value by a factor of σ . This is done at each iteration, so to satisfy

$$p_k = \min(p_{des}, \sigma \cdot p_{k-1}) \quad (14)$$

where p_{des} corresponds to the desired l_p norm. The implementation of each iteration follows Karlovitz's method using the particular selection of p given by (14).

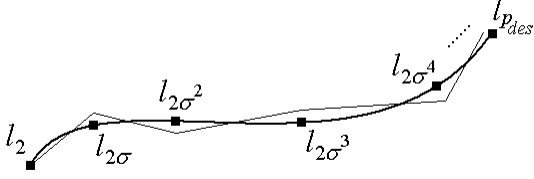


Fig. 1. Homotopy approach for IRLS l_p filter design.

It is worth noting that the method outlined above combines several ideas into a powerful approach. By not solving the desired l_p problem from the first iteration, one avoids the potential issues of Newton-based methods where convergence is guaranteed within a radius of convergence. It is well known that for $2 \leq p \leq \infty$ there exists a continuum of l_p solutions (as shown in Figure 1). By slowly increasing p from iteration to iteration one hopes to follow the continuum of solutions from l_2 towards the desired p . By choosing a reasonable σ the method can only spend one iteration at any given p and still remain close enough to the optimal path. Once the algorithm reaches a neighborhood of the desired p , it can be allowed to iterate at such p , in order to converge to the optimal solution. This process is analogous to *homotopy*, a commonly used family of optimization methods [14].

While l_2 and l_∞ designs offer meaningful approaches to filter design, the Constrained Least Squares (CLS) problem offers an interesting tradeoff to both approaches [15]. In the context of filter design, the CLS problem seems to be first presented by John Adams [16] in 1991. The problem Adams posed is a *Quadratic Programming* (QP) problem, well suited for off-the-shelf QP tools like those based on Lagrange multiplier theory [16]. However, Adams posed the problem in such a way that a transition band is required. Burrus et al. presented a formulation [17]–[19] where only a *transition frequency* is required; the transition band is *induced*; it does indeed exist but is not specified (it adjusts itself optimally according to the constraint specifications). The method by Burrus et al. is based on Lagrange multipliers and the Karush-Kuhn-Tucker (KKT) conditions.

An alternative to the KKT-based method mentioned above is the use of IRLS methods where a suitable weighting function serves as the constraining function over frequencies that exceed the constraint tolerance. Otherwise no weights are used, effectively forcing a least-squares solution. While this idea has been suggested by Burrus et al., one of the main contributions of this work is a thorough investigation of this approach, as well as proper documentation of numerical results, theoretical findings and proper code.

E. Infinite Impulse Response (IIR) l_p design

In contrast to FIR filters, an *Infinite Impulse Response* (IIR) filter is defined by two ordered vectors $\mathbf{a} \in \mathbb{R}^N$ and $\mathbf{b} \in \mathbb{R}^{M+1}$ (where $0 < M, N < \infty$), with frequency response given by

$$H(\omega) = \frac{B(\omega)}{A(\omega)} = \frac{\sum_{n=0}^M b_n e^{-j\omega n}}{1 + \sum_{n=1}^N a_n e^{-j\omega n}}$$

Hence the general l_p approximation problem is

$$\min_{\mathbf{a}_n, \mathbf{b}_n} \left\| \frac{\sum_{n=0}^M b_n e^{-j\omega n}}{1 + \sum_{n=1}^N a_n e^{-j\omega n}} - D(\omega) \right\|_p \quad (15)$$

which can be posed as a weighted least squares problem of the form

$$\min_{\mathbf{a}_n, \mathbf{b}_n} \left\| w(\omega) \cdot \left(\frac{\sum_{n=0}^M b_n e^{-j\omega n}}{1 + \sum_{n=1}^N a_n e^{-j\omega n}} - D(\omega) \right) \right\|_2^2 \quad (16)$$

It is possible to design similar problems to the ones outlined in Section I-D for FIR filters. However, it is worth keeping in mind the additional complications that IIR design involves, including the nonlinear least squares problem presented in Section I-E1 below.

1) *Least squares IIR literature review*: The weighted nonlinear formulation presented in (16) suggests the possibility of taking advantage of the flexibilities in design from the FIR problems. However this point comes at the expense of having to solve at each iteration a weighted nonlinear l_2 problem. Solving least squares approximations with rational functions is a nontrivial problem that has been studied extensively in diverse areas including statistics, applied mathematics and electrical engineering. One of the contributions of this document is a presentation in Section III-B on the subject of l_2 IIR filter design that captures and organizes previous relevant work. It also sets the framework for the proposed methods used in this document.

In the context of IIR digital filters there are three main groups of approaches to (16). Section III-B1 presents relevant work in the form of traditional optimization techniques. These are methods derived mainly from the applied mathematics community and are in general efficient and well understood. However the generality of such methods occasionally comes at the expense of being inefficient for some particular problems. Among the methods found in literature, the Davidon-Fletcher-Powell (DFP) algorithm [20], the damped Gauss-Newton method [21], [22], the Levenberg-Marquardt algorithm [23], [24], and the method of Kumaresan [25], [26] form the basis of a number of methods to solve (15).

A different approach to (15) from traditional optimization methods consists in *linearizing* (16) by transforming the problem into a simpler, linear form. While in principle this proposition seems inadequate (as the original problem is being transformed), Section III-B2 presents some logical attempts at linearizing (16) and how they connect with the original problem. The concept of *equation error* (a **weighted** form of the *solution error* that one is actually interested in solving) has been introduced and employed by a number of authors. In the context of filter design, E. Levy [27] presented an equation error linearization formulation in 1959 applied to analog filters. An alternative equation error approach presented by C. S. Burrus [28] in 1987 is based on the methods by Prony [29] and Pade [30]. The method by Burrus can be applied to frequency domain digital filter design, and is used in selected stages in some of the algorithms presented in this work.

An extension of the equation error methods is the group of *iterative prefiltering* algorithms presented in Section III-B8. These methods build on equation error methods by weighting (or *prefiltering*) their equation error formulation iteratively, with the intention to converge to the minimum of the solution error. Sanathanan and Koerner [31] presented in 1963 an algorithm (SK) that builds on an extension of Levy's method by iterating on Levy's formulation. Sid-Ahmed, Chottera and Jullien [32] presented in 1978 a similar algorithm to the SK method but applied to the digital filter problem.

A popular and well understood method is the one by Steiglitz and McBride [33], [34] introduced in 1966. The SMB method is time-domain based, and has been extended to a number of applications, including the frequency domain filter design problem [35]. Steiglitz and McBride used a two-phase method based on linearization. Initially (in *Mode-1*) their algorithm is essentially that of Sanathanan and Koerner but in time. This approach often diverges when close to

the solution; therefore their method can optionally switch to *Mode-2*, where a more traditional derivative-based approach is used.

A more recent linearization algorithm was presented by L. Jackson [36] in 2008. His approach is an iterative prefiltering method based directly in frequency domain, and uses diagonalization of certain matrices for efficiency.

While intuitive and relatively efficient, most linearization methods share a common problem: they often diverge close to the solution (this effect has been noted by a number of authors; a thorough review is presented in [35]). Section III-B13 presents the *quasilinearization* method derived by A. Soewito [35] in 1990. This algorithm is robust, efficient and well-tailored for the least squares IIR problem, and is the method of choice for this work.

II. FINITE IMPULSE RESPONSE FILTERS

This chapter discusses the problem of designing Finite Impulse Response (FIR) digital filters according to the l_p error criterion using Iterative Reweighted Least Squares methods. Section II-A gives an introduction to FIR filter design, including an overview of traditional FIR design methods. For the purposes of this work we are particularly interested in l_2 and l_∞ design methods, and their relation to relevant l_p design problems. Section II-B formally introduces the linear phase problem and presents results that are common to most of the problems considered in this work. Finally, Sections II-C through II-E present the application of the Iterative Reweighted Least Squares algorithm to other important problems in FIR digital filter design, including the relevant contributions of this work.

A. Traditional design of FIR filters

Section I-A introduced the notion of digital filters and filter design. In a general sense, an FIR filter design problem has the form

$$\min_{\mathbf{h}} \|f(\mathbf{h})\|$$

where $f(\cdot)$ defines an error function that depends on \mathbf{h} , and $\|\cdot\|$ is an arbitrary norm. While one could come up with a number of error formulations for digital filters, this chapter elaborates on the most commonly used, namely the linear phase and complex problems (both satisfy the linear form $f(\mathbf{h}) = \mathbf{D} - \mathbf{C}\mathbf{h}$ as will be shown later in this chapter). As far as norms, typically the l_2 and l_∞ norms are used. One of the contributions of this work is to demonstrate the usefulness of the more general l_p norms and their feasibility by using efficient IRLS-based algorithms.

1) *Traditional design of least squares (l_2) FIR filters:* Typically, FIR filters are designed by *discretizing* a desired frequency response $H_d(\omega)$ by taking L frequency samples at $\{\omega_0, \omega_1, \dots, \omega_{L-1}\}$. One could simply take the inverse Fourier transform of these samples and obtain L filter coefficients; this approach is known as the *Frequency Sampling design method* [28], which basically interpolates the frequency spectrum over the samples. However, it is often more desirable to take a large number of samples to design a small filter (large in the sense that $L \gg N$, where L is the number of frequency samples and N is the filter order). The weighted least-squares (l_2) norm (which considers the error energy) is defined by

$$\varepsilon_2 \triangleq \|\epsilon(\omega)\|_2 = \left(\frac{1}{\pi} \int_0^\pi W(\omega) |D(\omega) - H(\omega)|^2 d\omega \right)^{\frac{1}{2}} \quad (17)$$

where $D(\omega)$ and $H(\omega) = \mathcal{F}(\mathbf{h})$ are the desired and designed amplitude responses respectively. By acknowledging the convexity of (17), one can drop the root term; therefore a discretized form of (17) is given by

$$\varepsilon_2 = \sum_{k=0}^{L-1} W(\omega_k) |D(\omega_k) - H(\omega_k)|^2 \quad (18)$$

The solution of Equation (18) is given by

$$\mathbf{h} = \left(\mathbf{C}^T \mathbf{W}^T \mathbf{W} \mathbf{C} \right)^{-1} \mathbf{C}^T \mathbf{W}^T \mathbf{W} \mathbf{D} \quad (19)$$

where $\mathbf{W} = \text{diag}(\sqrt{w})$ contains the weighting vector w . By solving (19) one obtains an optimal l_2 approximation to the desired frequency response $D(\omega)$. Further discussion and other variations on least squares FIR design can be found in [28].

2) *Traditional design of minimax (l_∞) FIR filters:* In contrast to l_2 design, an l_∞ filter minimizes the maximum error across the designed filter's frequency response. A formal formulation of the problem [37], [38] is given by

$$\min_{\mathbf{h}} \max_{\omega} |D(\omega) - H(\omega; \mathbf{h})| \quad (20)$$

A discrete version of (20) is given by

$$\min_{\mathbf{h}} \max_k |D(\omega_k) - C_k \mathbf{h}| \quad (21)$$

Within the scope of filter design, the most commonly approach to solving (21) is the use of the *Alternation Theorem* [39], in the context of linear phase filters (to be discussed in Section II-B). In a nutshell the alternation theorem states that for a length- N FIR linear phase filter there are at least $N + 1$ *extrema points* (or frequencies). The Remez exchange algorithm [28], [37], [38] aims at finding these extrema frequencies iteratively, and is the most commonly used method for the minimax linear phase FIR design problem. Other approaches use more standard linear programming methods including the Simplex algorithm [40], [41] or interior point methods such as Karmarkar's algorithm [42].

The l_∞ problem is fundamental in filter design. While this document is not aimed covering the l_∞ problem in depth, portions of this work are devoted to the use of IRLS methods for standard problems as well as some innovative uses of minimax optimization.

B. Linear phase l_p filter design

Linear phase FIR filters are important tools in signal processing. As will be shown below, they do not require the user to specify a phase response in their design (since the assumption is that the desired phase response is indeed linear). Besides, they satisfy a number of symmetry properties that allow for the reduction of dimensions in the optimization process, making them easier to design computationally. Finally, there are applications where a linear phase is desired as such behavior is more physically meaningful.

1) *Four types of linear phase filters:* The frequency response of an FIR filter $h(n)$ is given by

$$H(\omega) = \sum_{n=0}^{N-1} h(n) e^{-j\omega n}$$

In general, $H(\omega) = R(\omega) + jI(\omega)$ is a periodic complex function of ω (with period 2π). Therefore it can be written as follows,

$$\begin{aligned} H(\omega) &= R(\omega) + jI(\omega) \\ &= A(\omega) e^{j\phi(\omega)} \end{aligned} \quad (22)$$

where the *magnitude* response is given by

$$A(\omega) = |H(\omega)| = \sqrt{R(\omega)^2 + I(\omega)^2} \quad (23)$$

and the *phase* response is

$$\phi(\omega) = \arctan \left(\frac{I(\omega)}{R(\omega)} \right)$$

However $A(\omega)$ is not analytic and $\phi(\omega)$ is not continuous. From a computational point of view (22) would have better properties if

both $A(\omega)$ and $\phi(\omega)$ were continuous analytic functions of ω ; an important class of filters for which this is true is the class of *linear phase* filters [28].

Linear phase filters have a frequency response of the form

$$H(\omega) = A(\omega)e^{j\phi(\omega)} \quad (24)$$

where $A(\omega)$ is the real, continuous *amplitude* response of $H(\omega)$ and

$$\phi(\omega) = K_1 + K_2\omega$$

is a **linear** phase function in ω (hence the name); K_1 and K_2 are constants. The jumps in the phase response correspond to sign reversals in the magnitude resulting as defined in (23).

Consider a length- N FIR filter (assume for the time being that N is odd). Its frequency response is given by

$$\begin{aligned} H(\omega) &= \sum_{n=0}^{N-1} h(n)e^{-j\omega n} \\ &= e^{-j\omega M} \sum_{n=0}^{2M} h(n)e^{j\omega(M-n)} \end{aligned} \quad (25)$$

where $M = \frac{N-1}{2}$. Equation (25) can be written as follows,

$$\begin{aligned} H(\omega) &= e^{-j\omega M} [h(0)e^{j\omega M} + \dots + h(M-1)e^{j\omega} + h(M) \\ &\quad + h(M+1)e^{-j\omega} + \dots + h(2M)e^{-j\omega M}] \end{aligned} \quad (26)$$

It is clear that for an odd-length FIR filter to have the linear phase form described in (24), the term inside braces in (26) must be a real function (thus becoming $A(\omega)$). By imposing even symmetry on the filter coefficients about the midpoint ($n = M$), that is

$$h(k) = h(2M - k)$$

equation (26) becomes

$$H(\omega) = e^{-j\omega M} \left[h(M) + 2 \sum_{n=0}^{M-1} h(n) \cos \omega(M-n) \right] \quad (27)$$

Similarly, with odd symmetry (i.e. $h(k) = h(2M - k)$) equation (26) becomes

$$H(\omega) = e^{j(\frac{\pi}{2} - \omega M)} 2 \sum_{n=0}^{M-1} h(n) \sin \omega(M-n) \quad (28)$$

Note that the term $h(M)$ disappears as the symmetry condition requires that

$$h(M) = h(N - M - 1) = -h(M) = 0$$

Similar expressions can be obtained for an even-length FIR filter,

$$\begin{aligned} H(\omega) &= \sum_{n=0}^{N-1} h(n)e^{-j\omega n} \\ &= e^{-j\omega M} \sum_{n=0}^{\frac{N}{2}-1} h(n)e^{j\omega(M-n)} \end{aligned} \quad (29)$$

It is clear that depending on the combinations of N and the symmetry of $h(n)$, it is possible to obtain four types of filters [28], [43], [44]. Table I shows the four possible linear phase FIR filters described by (24), where the second column refers to the type of filter symmetry.

N Odd	Even	$A(\omega) = h(M) + 2 \sum_{n=0}^{M-1} h(n) \cos \omega(M-n)$ $\phi(\omega) = -\omega M$
	Odd	$A(\omega) = 2 \sum_{n=0}^{M-1} h(n) \sin \omega(M-n)$ $\phi(\omega) = \frac{\pi}{2} - \omega M$
N Even	Even	$A(\omega) = h(M) + 2 \sum_{n=0}^{\frac{N}{2}-1} h(n) \cos \omega(M-n)$ $\phi(\omega) = -\omega M$
	Odd	$A(\omega) = 2 \sum_{n=0}^{\frac{N}{2}-1} h(n) \sin \omega(M-n)$ $\phi(\omega) = \frac{\pi}{2} - \omega M$

TABLE I
THE FOUR TYPES OF LINEAR PHASE FIR FILTERS.

2) *IRLS-based methods*: Section II-B1 introduced linear phase filters in detail. In this section we cover the use of IRLS methods to design linear phase FIR filters according to the l_p optimality criterion. Recall from Section II-B1 that for any of the four types of linear phase filters their frequency response can be expressed as

$$H(\omega) = A(\omega)e^{j(K_1 + K_2\omega)}$$

Since $A(\omega)$ is a real continuous function as defined by Table I, one can write the linear phase l_p design problem as follows

$$\min_{\mathbf{a}} \|D(\omega) - A(\omega; \mathbf{a})\|_p^p \quad (30)$$

where \mathbf{a} relates to \mathbf{h} by considering the symmetry properties outlined in Table I. Note that the two objects from the objective function inside the l_p norm are real. By sampling (30) one can write the design problem as follows

$$\begin{aligned} \min_{\mathbf{a}} \sum_k |D(\omega_k) - A(\omega_k; \mathbf{a})|^p \\ \text{or} \\ \min_{\mathbf{a}} \sum_k |D_k - \mathbf{C}_k \mathbf{a}|^p \end{aligned} \quad (31)$$

where D_k is the k -th element of the vector \mathbf{D} representing the sampled desired frequency response $D(\omega_k)$, and \mathbf{C}_k is the k -th row of the trigonometric kernel matrix as defined by Table I.

One can apply the basic IRLS approach described in Section I-C to solve (31) by posing this problem as a weighted least squares one:

$$\min_{\mathbf{a}} \sum_k w_k |\mathbf{D}_k - \mathbf{C}_k \mathbf{a}|^2 \quad (32)$$

The main issue becomes iteratively finding suitable weights \mathbf{w} for (32) so that the algorithm converges to the optimal solution \mathbf{a}^* of the l_p problem (30). Existence of adequate weights is guaranteed by Theorem 1 as presented in Section I-C; finding these optimal weights is indeed the difficult part. Clearly a reasonable choice for \mathbf{w} is that which turns (32) into (31), namely

$$\mathbf{w} = |\mathbf{D} - \mathbf{C}\mathbf{a}|^{p-2}$$

Therefore the basic IRLS algorithm for problem (31) would be:

- 1) Initialize the weights \mathbf{w}_0 (a reasonable choice is to make them all equal to one).
- 2) At the i -th iteration the solution is given by

$$\mathbf{a}_{i+1} = [\mathbf{C}^T \mathbf{W}_i^T \mathbf{W}_i \mathbf{C}]^{-1} \mathbf{C}^T \mathbf{W}_i^T \mathbf{W}_i \mathbf{D} \quad (33)$$

3) Update the weights with

$$\mathbf{w}_{i+1} = |\mathbf{D} - \mathbf{C}\mathbf{a}_{i+1}|^{p-2}$$

4) Repeat the last steps until convergence is reached.

It is important to note that $\mathbf{W}_i = \text{diag}(\sqrt{\mathbf{w}_i})$. In practice it has been found that this approach has practical deficiencies, since the inversion required by (33) often leads to an ill-posed problem and, in most cases, convergence is not achieved.

As mentioned before, the basic IRLS method has drawbacks that make it unsuitable for practical implementations. Charles Lawson considered a version of this algorithm applied to the solution of l_∞ problems (for details refer to [4]). His method has linear convergence and is prone to problems with proportionately small residuals that could lead to zero weights and the need for restarting the algorithm. In the context of l_p optimization, Rice and Usow [5] built upon Lawson's method by adapting it to l_p problems. Like Lawson's methods, the algorithm by Rice and Usow updates the weights in a multiplicative manner; their method shares similar drawbacks with Lawson's. Rice and Usow defined

$$w_{i+1}(\omega) = w_i^\alpha(\omega) |\epsilon_i(\omega)|^\beta$$

where

$$\alpha = \frac{\gamma(p-2)}{\gamma(p-2)+1}$$

and

$$\beta = \frac{\alpha}{2\gamma} = \frac{p-2}{2\gamma(p-2)+2}$$

and follow the basic algorithm.

L. A. Karlovitz realized the computational problems associated with the basic IRLS method and improved on it by partially updating the filter coefficient vector. He defines

$$\hat{\mathbf{a}}_{i+1} = [\mathbf{C}^T \mathbf{W}_i^T \mathbf{W}_i \mathbf{C}]^{-1} \mathbf{C}^T \mathbf{W}_i^T \mathbf{W}_i \mathbf{D} \quad (34)$$

and uses $\hat{\mathbf{a}}$ in

$$\mathbf{a}_{i+1} = \lambda \hat{\mathbf{a}}_{i+1} + (1-\lambda) \mathbf{a}_i \quad (35)$$

where $\lambda \in [0, 1]$ is a partial step parameter that must be adjusted at each iteration. Karlovitz's method [8] has been shown to converge globally for even values of p (where $2 \leq p < \infty$). In practice, convergence problems have been found even under such assumptions. Karlovitz proposed the use of line searches to find the *optimal* value of λ at each iteration, which basically creates an independent optimization problem nested inside each iteration of the IRLS algorithm. While computationally this search process for the optimal λ makes Karlovitz's method impractical, his work indicates the feasibility of IRLS methods and proves that partial updating indeed overcomes some of the problems in the basic IRLS method. Furthermore, Karlovitz's method is the first one to depart from a multiplicative updating of the weights in favor of an additive updating on the filter coefficients. In this way some of the problems in the Lawson-Rice-Usow approach are overcome, especially the need for restarting the algorithm.

S. W. Kahng built upon the findings by Karlovitz by considering the process of finding an adequate λ for partial updating. He applied Newton-Raphson's method to this problem and proposed a closed form solution for λ , given by

$$\lambda = \frac{1}{p-1} \quad (36)$$

resulting in

$$\mathbf{a}_{i+1} = \lambda \hat{\mathbf{a}}_{i+1} + (1-\lambda) \mathbf{a}_i \quad (37)$$

The rest of Kahng's algorithm follows Karlovitz's approach. However, since λ is fixed, there is no need to perform the linear search

at each iteration. Kahng's method has an added benefit: since it uses Newton's method to find λ , the algorithm tends to converge much faster than previous approaches. It has indeed been shown to converge quadratically. However, Newton-Raphson-based algorithms are not guaranteed to converge globally unless at some point the existing solution lies close enough to the solution, within their radius of convergence [11]. Fletcher, Grant and Hebden [10] derived the same results independently.

Burrus, Barreto and Selesnick [7], [12], [13] modified Kahng's methods in several important ways in order to improve on their initial and final convergence rates and the method's stability (we refer to this method as BBS). The first improvement is analogous to a *homotopy* [14]. Up to this point all efforts in l_p filter design attempted to solve the *actual* l_p problem from the first iteration. In general there is no reason to believe that an initial guess derived from an unweighted l_2 formulation (that is, the l_2 design that one would get by setting $\mathbf{w}_0 = \hat{1}$) will look in any way similar to the actual l_p solution that one is interested in. However it is known that there exists a continuity of l_p solutions for $1 < p < \infty$. In other words, if \mathbf{a}_2^* is the optimal l_2 solution, there exists a p for which the optimal l_p solution \mathbf{a}_p^* is arbitrarily close to \mathbf{a}_2^* ; that is, for a given $\delta > 0$

$$\|\mathbf{a}_2^* - \mathbf{a}_p^*\| \leq \delta \quad \text{for some } p \in (2, \infty)$$

This fact allows anyone to *gradually move* from an l_p solution to an l_q solution.

To accelerate initial convergence, the BBS method of Burrus et al. initially solves for l_2 by setting $p_0 = 2$ and then sets $p_i = \sigma \cdot p_{i-1}$, where σ is a convergence parameter defined by $1 \leq \sigma \leq 2$. Therefore at the i -th iteration

$$p_i = \min(p_{des}, \sigma p_{i-1}) \quad (38)$$

where p_{des} corresponds to the desired l_p solution. The implementation of each iteration follows Karlovitz's method with Kahng's choice of λ , using the particular selection of p given by (38).

To summarize, define the class of IRLS algorithms as follows: after i iterations, given a vector \mathbf{a}_i the IRLS iteration requires two steps,

- 1) Find $\mathbf{w}_i = f(\mathbf{a}_i)$
- 2) Find $\mathbf{a}_{i+1} = g(\mathbf{w}_i, \mathbf{a}_i)$

The following is a summary of the IRLS-based algorithms discussed so far and their corresponding updating functions:

- 1) Basic IRLS algorithm.

- $\mathbf{w}_i = |\mathbf{D} - \mathbf{C}\mathbf{a}_i|^{p-2}$
- $\mathbf{W}_i = \text{diag}(\sqrt{\mathbf{w}_i})$
- $\mathbf{a}_{i+1} = [\mathbf{C}^T \mathbf{W}_i^T \mathbf{W}_i \mathbf{C}]^{-1} \mathbf{C}^T \mathbf{W}_i^T \mathbf{W}_i \mathbf{D}$

- 2) Rice-Usow-Lawson (RUL) method

- $\mathbf{w}_i = \mathbf{w}_{i-1}^\alpha |\mathbf{D} - \mathbf{C}\mathbf{a}_i|^{\frac{\alpha}{2\gamma}}$
- $\mathbf{W}_i = \text{diag}(\mathbf{w}_i)$
- $\mathbf{a}_{i+1} = [\mathbf{C}^T \mathbf{W}_i^T \mathbf{W}_i \mathbf{C}]^{-1} \mathbf{C}^T \mathbf{W}_i^T \mathbf{W}_i \mathbf{D}$
- $\alpha = \frac{\gamma(p-2)}{\gamma(p-2)+1}$
- γ constant

- 3) Karlovitz' method

- $\mathbf{w}_i = |\mathbf{D} - \mathbf{C}\mathbf{a}_i|^{p-2}$
- $\mathbf{W}_i = \text{diag}(\sqrt{\mathbf{w}_i})$
- $\mathbf{a}_{i+1} = \lambda [\mathbf{C}^T \mathbf{W}_i^T \mathbf{W}_i \mathbf{C}]^{-1} \mathbf{C}^T \mathbf{W}_i^T \mathbf{W}_i \mathbf{D} + (1-\lambda) \mathbf{a}_i$
- λ constant

- 4) Kahng's method

- $\mathbf{w}_i = |\mathbf{D} - \mathbf{C}\mathbf{a}_i|^{p-2}$
- $\mathbf{W}_i = \text{diag}(\sqrt{\mathbf{w}_i})$
- $\mathbf{a}_{i+1} = \left(\frac{1}{p-1}\right) [\mathbf{C}^T \mathbf{W}_i^T \mathbf{W}_i \mathbf{C}]^{-1} \mathbf{C}^T \mathbf{W}_i^T \mathbf{W}_i \mathbf{D} + \left(\frac{p-2}{p-1}\right) \mathbf{a}_i$

5) BBS method

- $p_i = \min(p_{des}, \sigma \cdot p_{i-1})$
- $w_i = |D - Ca_i|^{p_i-2}$
- $W_i = \text{diag}(\sqrt{w_i})$
- $a_{i+1} = \left(\frac{1}{p_i-1} \right) [C^T W_i^T W_i C]^{-1} C^T W_i^T W_i D + \left(\frac{p_i-2}{p_i-1} \right) a_i$
- σ constant

3) *Modified adaptive IRLS algorithm*: Much of the performance of a method is based upon whether it can actually converge given a certain error measure. In the case of the methods described above, both convergence rate and stability play an important role in their performance. Both Karlovitz and RUL methods are supposed to converge linearly, while Kahng's and the BBS methods converge quadratically, since they both use a Newton-based additive update of the weights.

Barreto showed in [12] that the modified version of Kahng's method (or BBS) typically converges faster than the RUL algorithm. However, this approach presents some peculiar problems that depend on the transition bandwidth β . For some particular values of β , the BBS method will result in an ill-posed weight matrix that causes the l_p error to increase dramatically after a few iterations as illustrated in Figure 2 (where $f = \omega/2\pi$).

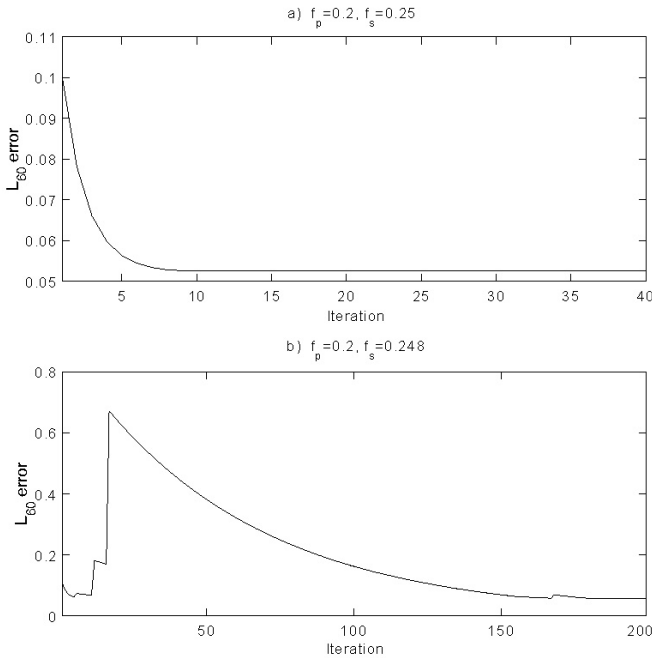


Fig. 2. Error jumps on IRLS methods.

Two facts can be derived from the examples in Figure 2: for this particular bandwidth the error increased slightly after the fifth and eleventh iterations, and increased dramatically after the sixteenth. Also, it is worth to notice that after such increase, the error started to decrease quadratically and that, at a certain point, the error became flat (thus reaching the numerical accuracy limits of the digital system).

The effects of different values of σ were studied to find out if a relationship between σ and the error increase could be determined. Figure 3 shows the l_p error for different values of β and for $\sigma = 1.7$. It can be seen that some particular bandwidths cause the algorithm to produce a very large error.

Our studies (as well as previous work from J. A. Barreto [12]) demonstrate that this error explosion occurs only for a small range of

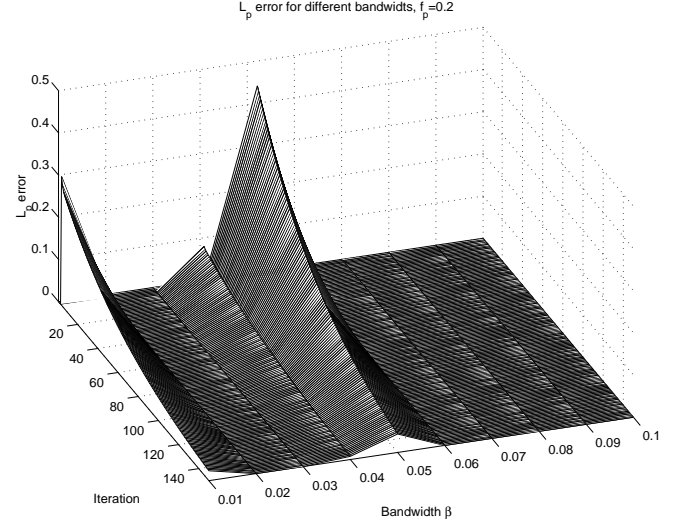


Fig. 3. Relationship between bandwidth and error jumps.

bandwidth specifications. Under most circumstances the BBS method exhibits fast convergence properties to the desired solution. However at this point it is not understood what causes the error increase and therefore this event cannot be anticipated. In order to avoid such problem, we propose the use of an adaptive scheme that modifies the BBS step. As p increases the step from a current l_p guess to the next also increases, as described in (38). In other words, at the i -th iteration one approximates the $l_{2\sigma^i}$ solution (as long as the algorithm has not yet reached the desired p); the next iteration one approximates $l_{2\sigma^{i+1}}$. There is always a possibility that these two solutions lie far apart enough that the algorithm takes a descent step so that the $l_{2\sigma^{i+1}}$ guess is too far away from the actual $l_{2\sigma^{i+1}}$ solution. This is better illustrated in Figure 4.

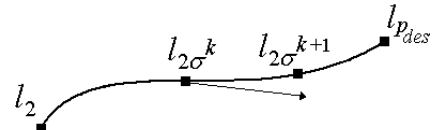


Fig. 4. A step too long for IRLS methods.

The conclusions derived above suggest the possibility to use an adaptive algorithm [45] that changes the value of σ so that the error always decreases. This idea was implemented by calculating temporary new weight and filter coefficient vectors that will not become the updated versions unless their resulting error is smaller than the previous one. If this is not the case, the algorithm "tries" two values of σ , namely

$$\sigma_L = \sigma * (1 - \delta) \quad \text{and} \quad \sigma_H = \sigma * (1 + \delta) \quad (39)$$

(where δ is an updating variable). The resulting errors for each attempt are calculated, and σ is updated according to the value that produced the smallest error. The error of this new σ is compared to the error of the nonupdated weights and coefficients, and if the new σ produces a smaller error, then such vectors are updated; otherwise another update of σ is performed. The *modified adaptive IRLS algorithm* can be summarized as follows,

- 1) Find the unweighted approximation $a_0 = [C^T C]^{-1} C^T D$ and use $p_0 = 2\sigma$ (with $1 \leq \sigma \leq 2$)

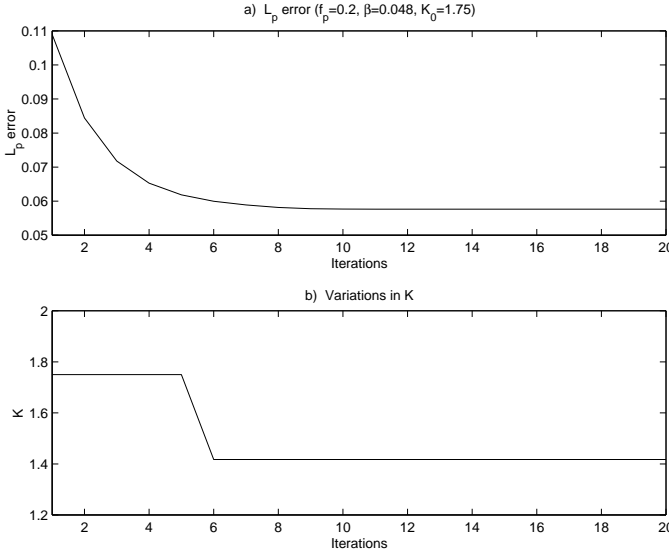


Fig. 5. FIR design example using adaptive method. a) l_p error obtained with the adaptive method; b) Change of σ .

- 2) Iteratively solve (34) and (35) using $\lambda_i = \frac{1}{p_i-1}$ and find the resulting error ε_i for the i -th iteration
- 3) If $\varepsilon_i \gg \varepsilon_{i-1}$,
 - Calculate (39)
 - Select the smallest of ε_{σ_L} and ε_{σ_H} to compare it with ε_i until a value is found that results in a decreasing error

Otherwise iterate as in the BBS algorithm.

The algorithm described above changes the value of σ that causes the algorithm to produce a large error. The value of σ is updated as many times as necessary without changing the values of the weights, the filter coefficients, or p . If an optimal value of σ exists, the algorithm will find it and continue with this new value until another update in σ becomes necessary.

The algorithm described above was implemented for several combinations of σ and β ; for all cases the new algorithm converged faster than the BBS algorithm (unless the values of σ and β are such that the error never increases). The results are shown in Figure 5.a for the specifications from Figure 2. Whereas using the BBS method for this particular case results in a large error after the sixteenth iteration, the adaptive method converged before ten iterations.

Figure 5.b illustrates the change of σ per iteration in the adaptive method, using an update factor of $\delta = 0.1$. The l_p error stops decreasing after the fifth iteration (where the BBS method introduces the large error); however, the adaptive algorithm adjusts the value of σ so that the l_p error continues decreasing. The algorithm decreased the initial value of σ from 1.75 to its final value of 1.4175 (at the expense of only one additional iteration with $\sigma = 1.575$).

One result worth noting is the relationship between l_2 and l_∞ solutions and how they compare to l_p designs. Figure 6 shows a comparison of designs for a length-21 Type-I linear phase low pass FIR filter with transition band defined by $f = \{0.2, 0.24\}$. The curve shows the l_2 versus l_∞ errors (namely ε_2 and ε_∞); the values of p used to make this curve were $p = \{2, 2.2, 2.5, 3, 4, 5, 7, 10, 15, 20, 30, 50, 60, 100, 150, 200, 400, \infty\}$ (Matlab's `firls` and `firpm` functions were used to design the l_2 and l_∞ filters respectively). Note the very small decrease in ε_∞ after p reaches 100. The curve suggests that a better compromise between ε_2 and ε_∞ can be reached by choosing $2 < p < \infty$. Furthermore, to get better results one can concentrate on values between $p = 5$

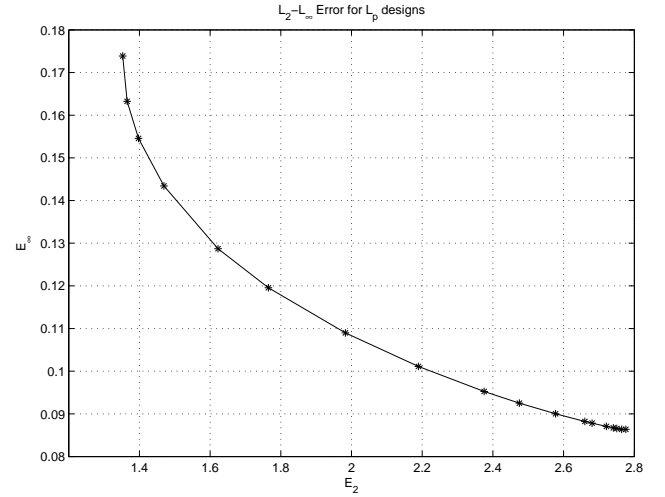


Fig. 6. Relationship between l_2 and l_∞ errors for l_p FIR filter design.

and $p = 20$; fortunately, for values of p so low no numerical complications arise and convergence is reached in a few iterations.

C. Complex l_p problem

The design of linear phase filters has been intensively discussed in literature. For the two most common error criteria (l_2 and l_∞), optimal solution algorithms exist. The least squares norm filter can be found by solving an overdetermined system of equations, whereas the Chebishev norm filter is easily found by using either the Remez algorithm or linear programming. For many typical applications, linear phase filters are good enough; however, when arbitrary magnitude and phase constraints are required, a more complicated approach must be taken since such design results in a complex approximation problem. By replacing \mathbf{C} in the linear phase algorithm with a complex Fourier kernel matrix, and the real desired frequency vector \mathbf{D} with a complex one, one can use the same algorithm from Section II-B3 to design complex l_p filters.

D. Magnitude l_p problem

In some applications, the effects of phase are not a necessary factor to consider when designing a filter. For these applications, control of the filter's magnitude response is a priority for the designer. In order to improve the magnitude response of a filter, one must not explicitly include a phase, so that the optimization algorithm can look for the best filter that approximates a specified magnitude, without being constrained about optimizing for a phase response too.

1) *Power approximation formulation:* The magnitude approximation problem can be formulated as follows:

$$\min_{\mathbf{h}} \|D(\omega) - |H(\omega; \mathbf{h})|\|_p^p \quad (40)$$

Unfortunately, the second term inside the norm (namely the absolute value function) is not differentiable when its argument is zero. Although one could propose ways to work around this problem, we propose the use of a different design criterion, namely the approximation of a desired magnitude squared. The resulting problem is

$$\min_{\mathbf{h}} \|D(\omega)^2 - |H(\omega; \mathbf{h})|^2\|_p^p$$

The autocorrelation $r(n)$ of a causal length- N FIR filter $h(n)$ is given by

$$r(n) = h(n) * h(-n) = \sum_{k=-(N-1)}^{N-1} h(k)h(n+k) \quad (41)$$

The Fourier transform of the autocorrelation $r(n)$ is known as the *Power Spectral Density* function [46] $R(\omega)$ (or simply the SPD), and is defined as follows,

$$\begin{aligned} R(\omega) &= \sum_{n=-(N-1)}^{N-1} r(n)e^{-j\omega n} \\ &= \sum_{n=-(N-1)}^{N-1} \sum_{k=-(N-1)}^{N-1} h(n)h(n+k)e^{-j\omega n} \end{aligned}$$

From the properties of the Fourier Transform [47, §3.3] one can show that there exists a frequency domain relationship between $h(n)$ and $r(n)$ given by

$$R(\omega) = H(\omega) \cdot H^*(-\omega) = |H(\omega)|^2$$

This relationship suggests a way to design magnitude-squared filters, namely by using the filter's autocorrelation coefficients instead of the filter coefficients themselves. In this way, one can avoid the use of the non-differentiable magnitude response.

An important property to note at this point is the fact that since the filter coefficients are real, one can see from (41) that the autocorrelation function $r(n)$ is symmetric; thus it is sufficient to consider its last N values. As a result, the PSD can be written as

$$R(\omega) = \sum_n r(n)e^{-j\omega n} = r(0) + \sum_{n=1}^{N-1} 2r(n) \cos \omega n$$

in a similar way to the linear phase problem.

The symmetry property introduced above allows for the use of the l_p linear phase algorithm of section (II-B) to obtain the autocorrelation coefficients of $h(n)$. However, there is an important step missing in this discussion: how to obtain the filter coefficients from its autocorrelation. To achieve this goal, one can follow a procedure known as *Spectral Factorization*. The objective is to use the autocorrelation coefficients $\mathbf{r} \in \mathbb{R}^N$ instead of the filter coefficients $\mathbf{h} \in \mathbb{R}^N$ as the optimization variables. The variable transformation is done using (42), which is not a one-to-one transformation. Because of the last result, there is a necessary condition for a vector $\mathbf{r} \in \mathbb{R}^N$ to be a valid autocorrelation vector of a filter. This is summarized [48] in the *spectral factorization theorem*, which states that $\mathbf{r} \in \mathbb{R}^N$ is the autocorrelation function of a filter $h(n)$ if and only if $R(\omega) \geq 0$ for all $\omega \in [0, \pi]$. This turns out to be a necessary and sufficient condition [48] for the existence of $r(n)$. Once the autocorrelation vector \mathbf{r} is found using existing robust interior-point algorithms, the filter coefficients can be calculated via spectral factorization techniques.

Assuming a valid vector $\mathbf{r} \in \mathbb{R}^N$ can be found for a particular filter \mathbf{h} , the problem presented in (40) can be rewritten as

$$L(\omega)^2 \leq R(\omega) \leq U(\omega)^2 \quad \forall \omega \in [0, \pi] \quad (42)$$

In (42) the existence condition $R(\omega) \geq 0$ is redundant since $0 \leq L(\omega)^2$ and, thus, is not included in the problem definition. For each ω , the constraints of (42) constitute a pair of linear inequalities in the vector \mathbf{r} ; therefore the constraint is convex in \mathbf{r} . Thus the change of variable transforms a nonconvex optimization problem in \mathbf{h} into a convex problem in \mathbf{r} .

E. l_p error as a function of frequency

Previous sections have discussed the importance of complex least-square and Chebishev error criteria in the context of filter design. In many applications any of these two approaches would provide adequate results. However, a case could be made where one might want to minimize the error energy in a range of frequencies while keeping control of the maximum error in a different band. This idea results particularly interesting when one considers the use of different l_p norms in different frequency bands. In principle one would be interested in solving

$$\min_{\mathbf{h}} \|D(\omega_{pb}) - H(\omega_{pb}; \mathbf{h})\|_p + \|D(\omega_{sb}) - H(\omega_{sb}; \mathbf{h})\|_q \quad (43)$$

where $\{\omega_{pb} \in \Omega_{pb}, \omega_{sb} \in \Omega_{sb}\}$ represent the pass and stopband frequencies respectively. In principle one would want $\Omega_{pb} \cap \Omega_{sb} = \{\emptyset\}$. Therefore problem (43) can be written as

$$\begin{aligned} \min_{\mathbf{h}} & \sqrt[p]{\sum_{\omega_{pb}} |D(\omega_{pb}) - H(\omega_{pb}; \mathbf{h})|^p} \\ & + \sqrt[q]{\sum_{\omega_{sb}} |D(\omega_{sb}) - H(\omega_{sb}; \mathbf{h})|^q} \end{aligned} \quad (44)$$

One major obstacle in (44) is the presence of the roots around the summation terms. These roots prevent us from writing (44) in a simple vector form. Instead, one can consider the use of a similar *metric* function as follows

$$\begin{aligned} \min_{\mathbf{h}} & \sum_{\omega_{pb}} |D(\omega_{pb}) - H(\omega_{pb}; \mathbf{h})|^p + \\ & \sum_{\omega_{sb}} |D(\omega_{sb}) - H(\omega_{sb}; \mathbf{h})|^q \end{aligned} \quad (45)$$

This expression is similar to (44) but does not include the root terms. An advantage of using the IRLS approach on (45) is that one can formulate this problem in the frequency domain and properly separate residual terms from different bands into different vectors. In this manner, the l_p modified measure given by (45) can be made into a frequency-dependent function of $p(\omega)$ as follows,

$$\min_{\mathbf{h}} \|D(\omega) - H(\omega; \mathbf{h})\|_{p(\omega)}^{p(\omega)} = \sum_{\omega} |D(\omega) - H(\omega; \mathbf{h})|^{p(\omega)}$$

Therefore this *frequency-varying* l_p problem can be solved following the modified IRLS algorithm outlined in Section II-B3 with the following modification: at the i -th iteration the weights are updated according to

$$\mathbf{w}_i = |\mathbf{D} - \mathbf{C}\mathbf{a}_i|^{p(\omega)-2}$$

It is fundamental to note that the proposed method does not indeed solve a linear combination of l_p norms. In fact, it can be shown that the expression (45) is not a norm but a metric. While from a theoretical perspective this fact might make (45) a less interesting distance, as it turns out one can use (45) to solve the far more interesting CLS problem, as discussed below in Section II-F.

F. Constrained Least Squares (CLS) problem

One of the common obstacles to innovation occurs when knowledge settles on a particular way of dealing with problems. While new ideas keep appearing suggesting innovative approaches to design digital filters, it is all too common in practice that l_2 and l_∞ dominate error criteria specifications. This section is devoted to exploring a different way of thinking about digital filters. It is important to note that up to this point we are not discussing an algorithm yet. The main concern being brought into play here is the specification (or description) of the design problem. Once the *Constrained Least*

Squares (CLS) problem formulation is introduced, we will present an IRLS implementation to solve it, and will justify our approach over other existing approaches. It is the author's belief that under general conditions one should always use our IRLS implementation over other methods, especially when considering the associated management of transition regions.

The CLS problem was introduced in Section I-D and is repeated here for clarity,

$$\begin{aligned} \min_{\mathbf{h}} \quad & \|D(\omega) - H(\omega; \mathbf{h})\|_2 \\ \text{subject to} \quad & |D(\omega) - H(\omega; \mathbf{h})| \leq \tau \end{aligned} \quad (46)$$

To the best of our knowledge this problem was first introduced in the context of filter design by John Adams [16] in 1991. The main idea consists in approximating iteratively a desired frequency response in a least squares sense except in the event that any frequency exhibits an error larger than a specified tolerance τ . At each iteration the problem is adjusted in order to reduce the error on offending frequencies (i.e. those which do not meet the constraint specifications). Ideally, convergence is reached when the *altered* least squares problem has a frequency response whose error does not exceed constraint specifications. As will be shown below, this goal might not be attained depending on how the problem is posed.

Adams and some collaborators have worked in this problem and several variations [15]. However his main (and original) problem was illustrated in [16] with the following important assumption: *the definition of a desired frequency response must include a fixed non-zero width transition band*. His method uses Lagrange multiplier theory and alternation methods to find frequencies that exceed constraints and minimize the error at such locations, with an overall least squares error criterion.

Burrus, Selesnick and Lang [17] looked at this problem from a similar perspective, but relaxed the design specifications so that only a *transition frequency* needs to be specified. The actual transition band does indeed exist, and it centers itself around the specified transition frequency; its width adjusts as the algorithm iterates (constraint tolerances are still specified). Their solution method is similar to Adams' approach, and explicitly uses the Karush-Kuhn-Tucker (KKT) conditions together with an alternation method to minimize the least squares error while constraining the maximum error to meet specifications.

C. S. Burrus and the author of this work have been working on the CLS problem using IRLS methods with positive results. This document is the first thorough presentation of the method, contributions, results and code for this approach, and constitutes one of the main contributions of this work. It is crucial to note that there are two separate issues in this problem: on one hand there is the matter of the actual problem formulation, mainly depending on whether a transition band is specified or not; on the other hand there is the question of how the selected problem description is actually met (what algorithm is used). Our approach follows the problem description by Burrus et al. shown in [17] with an IRLS implementation.

1) *Two problem formulations*: As mentioned in Section II-F, one can address problem (46) in two ways depending on how one views the role of the transition band in a CLS problem. The original problem posed by Adams in [16] can be written as follows,

$$\begin{aligned} \min_{\mathbf{h}} \quad & \|D(\omega) - H(\omega; \mathbf{h})\|_2 \\ \text{s.t.} \quad & |D(\omega) - H(\omega; \mathbf{h})| \leq \tau \quad \forall \omega \in [0, \omega_{pb}] \cup [\omega_{sb}, \pi] \end{aligned} \quad (47)$$

where $0 < \omega_{pb} < \omega_{sb} < \pi$. From a traditional standpoint this formulation feels familiar. It assigns **fixed frequencies** to the transition band edges as a number of filter design techniques do. As it turns out, however, one might not want to do this in CLS design.

An alternate formulation to (47) could implicitly introduce a *transition frequency* ω_{tb} (where $\omega_{pb} < \omega_{tb} < \omega_{sb}$); the user only specifies ω_{tb} . Consider

$$\begin{aligned} \min_{\mathbf{h}} \quad & \|D(\omega) - H(\omega; \mathbf{h})\|_2 \quad \forall \omega \in [0, \pi] \\ \text{subject to} \quad & |D(\omega) - H(\omega; \mathbf{h})| \leq \tau \quad \forall \omega \in [0, \omega_{pb}] \cup [\omega_{sb}, \pi] \end{aligned} \quad (48)$$

The algorithm at each iteration generates an *induced transition band* in order to satisfy the constraints in (48). Therefore $\{\omega_{pb}, \omega_{sb}\}$ vary at each iteration.

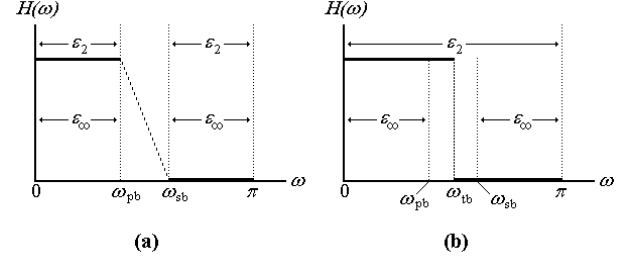


Fig. 7. Two formulations for Constrained Least Squares problems.

It is critical to point out the differences between (47) and (48). Figure 7.a explains Adams' CLS formulation, where the desired filter response is only specified at the fixed pass and stop bands. At any iteration, Adam's method attempts to minimize the least squares error (ϵ_2) at both bands while trying to satisfy the constraint τ . Note that one could think of the constraint requirements in terms of the Chebishev error ϵ_∞ by writing (47) as follows,

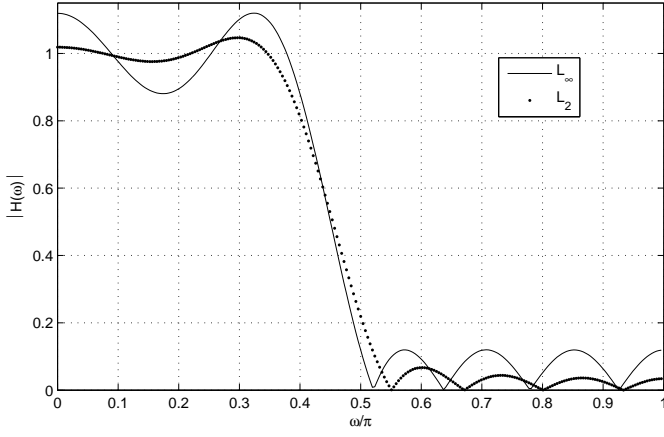
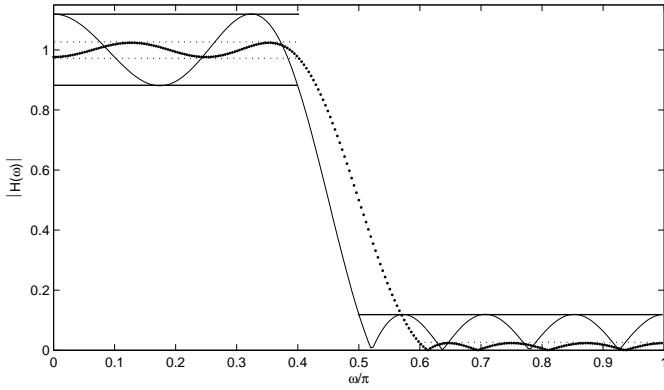
$$\begin{aligned} \min_{\mathbf{h}} \quad & \|D(\omega) - H(\omega; \mathbf{h})\|_2 \\ \text{s.t.} \quad & \|D(\omega) - H(\omega; \mathbf{h})\|_\infty \leq \tau \quad \forall \omega \in [0, \omega_{pb}] \cup [\omega_{sb}, \pi] \end{aligned}$$

In contrast, Figure 7.b illustrates our proposed problem (48). The idea is to minimize the least squared error ϵ_2 across **all** frequencies while ensuring that constraints are met in an intelligent manner. At this point one can think of the interval $(\omega_{pb}, \omega_{sb})$ as an *induced* transition band, useful for the purposes of constraining the filter. Section II-F2 presents the actual algorithms that solve (48), including the process of finding $\{\omega_{pb}, \omega_{sb}\}$.

It is important to note an interesting behavior of transition bands and extrema points in l_2 and l_∞ filters. Figure 8 shows l_2 and l_∞ length-15 linear phase filters (designed using Matlab's `firls` and `firpm` functions); the transition band was specified at $\{\omega_{pb} = 0.4/\pi, \omega_{sb} = 0.5/\pi\}$. The dotted l_2 filter illustrates an important behavior of least squares filters: typically the maximum error of an l_2 filter is located at the transition band. The solid l_∞ filter shows why minimax filters are important: despite their larger error across most of the bands, the filter shows the same maximum error at all extrema points, including the transition band edge frequencies. In a CLS problem then, typically an algorithm will attempt to reduce iteratively the maximum error (usually located around the transition band) of a series of least squares filters.

Another important fact results from the relationship between the transition band width and the resulting error amplitude in l_∞ filters. Figure 9 shows two l_∞ designs; the transition bands were set at $\{0.4/\pi, 0.5/\pi\}$ for the solid line design, and at $\{0.4/\pi, 0.6/\pi\}$ for the dotted line one. One can see that by widening the transition band a decrease in error ripple amplitude is induced.

These two results together illustrate the importance of the transition bandwidth for a CLS design. Clearly one can decrease maximum error tolerances by widening the transition band. Yet finding the perfect balance between a transition bandwidth and a given tolerance

Fig. 8. Comparison of l_2 and l_∞ filters.Fig. 9. Effects of transition bands in l_∞ filters.

can prove a difficult task, as will be shown in Section II-F2. Hence the relevance of a CLS method that is not restricted by two types of specifications competing against each other. In principle, one should just determine how much error one can live with, and allow an algorithm to find the optimal transition band that meets such tolerance.

2) *Two problem solutions:* Section II-F1 introduced some important remarks regarding the behavior of extrema points and transition bands in l_2 and l_∞ filters. As one increases the constraints on an l_2 filter, the result is a filter whose frequency response looks more and more like an l_∞ filter.

Section II-E introduced the frequency-varying problem and an IRLS-based method to solve it. It was also mentioned that, while the method does not solve the intended problem (but a similar one), it could prove to be useful for the CLS problem. As it turns out, in CLS design one is merely interested in solving an unweighted, constrained least squares problem. In this work, we achieve this by solving a sequence of weighted, unconstrained least squares problems, where the sole role of the weights is to "constraint" the maximum error of the frequency response at each iteration. In other words, one would like to find weights w such that

$$\begin{aligned} \min_{\mathbf{h}} \quad & \|D(\omega) - H(\omega; \mathbf{h})\|_2 \\ \text{s.t.} \quad & \|D(\omega) - H(\omega; \mathbf{h})\|_\infty \leq \tau \quad \forall \omega \in [0, \omega_{pb}] \cup [\omega_{sb}, \pi] \end{aligned}$$

is equivalent to

$$\min_{\mathbf{h}} \|w(\omega) \cdot (D(\omega) - H(\omega; \mathbf{h}))\|_2$$

Hence one can revisit the frequency-varying design method and use

it to solve the CLS problem. Assuming that one can reasonably approximate l_∞ by using high values of p , at each iteration the main idea is to use an l_p weighting function only at frequencies where the constraints are exceeded. A formal formulation of this statement is

$$w(\epsilon(\omega)) = \begin{cases} |\epsilon(\omega)|^{\frac{p-2}{2}} & \text{if } |\epsilon(\omega)| > \tau \\ 1 & \text{otherwise} \end{cases}$$

Assuming a suitable weighting function existed such that the specified tolerances are related to the frequency response constraints, the IRLS method would iterate and assign rather large weights to frequencies exceeding the constraints, while inactive frequencies get a weight of one. As the method iterates, frequencies with large errors move the response closer to the desired tolerance. Ideally, all the active constraint frequencies would eventually meet the constraints. Therefore the task becomes to find a suitable weighting function that *penalizes* large errors in order to have all the frequencies satisfying the constraints; once this condition is met, we have reached the desired solution.

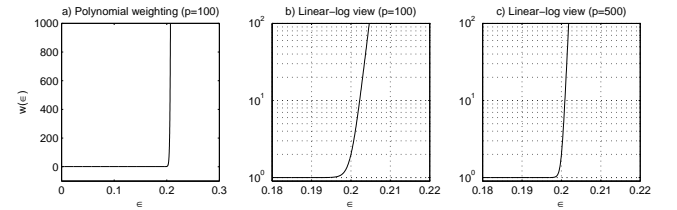
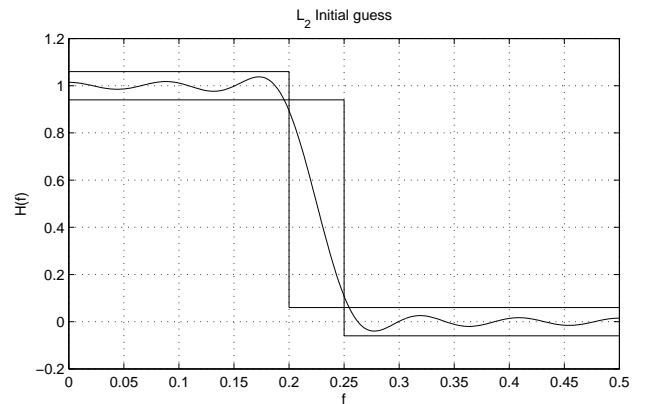


Fig. 10. CLS polynomial weighting function.

One proposed way to find adequate weights to meet constraints is given by a *polynomial weighting function* of the form

$$w(\omega) = 1 + \left| \frac{\epsilon(\omega)}{\tau} \right|^{\frac{p-2}{2}}$$

where τ effectively serves as a threshold to determine whether a weight is dominated by either unity or the familiar l_p weighting term. Figure 10 illustrates the behavior of such a curve.

Fig. 11. Original l_2 guess for CLS algorithm.

In practice the method outlined above has proven robust particularly in connection with the specified transition band design. Consider the least squares design in Figure 11 (using a length-21 Type-I linear phase low-pass FIR filter with linear transition frequencies $\{0.2, 0.25\}$). This example illustrates the typical effect of CLS methods over l_2 designs; the largest error (in an l_∞ sense) can be located at the edges of the transition band. Figures 12 and 13 illustrate

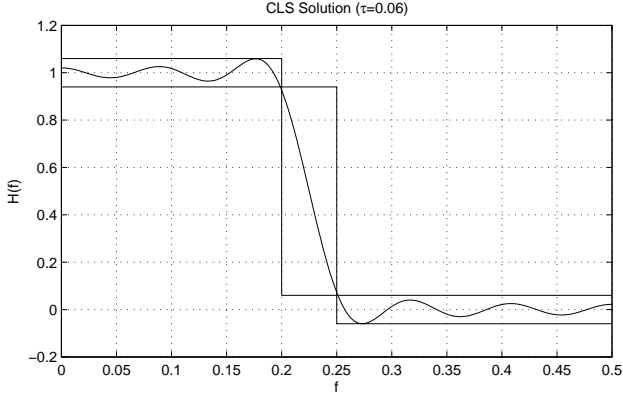


Fig. 12. CLS design example using mild constraints.

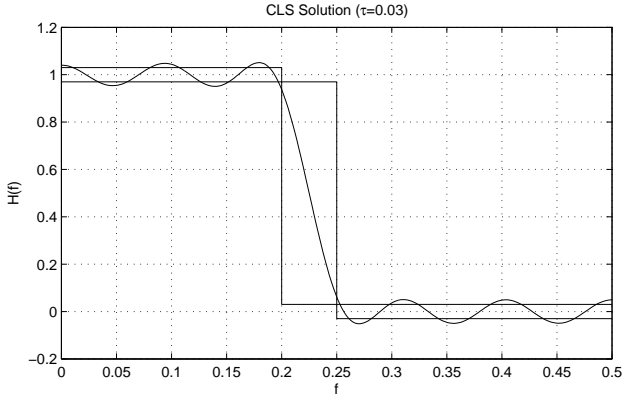


Fig. 13. CLS design example using tight constraints.

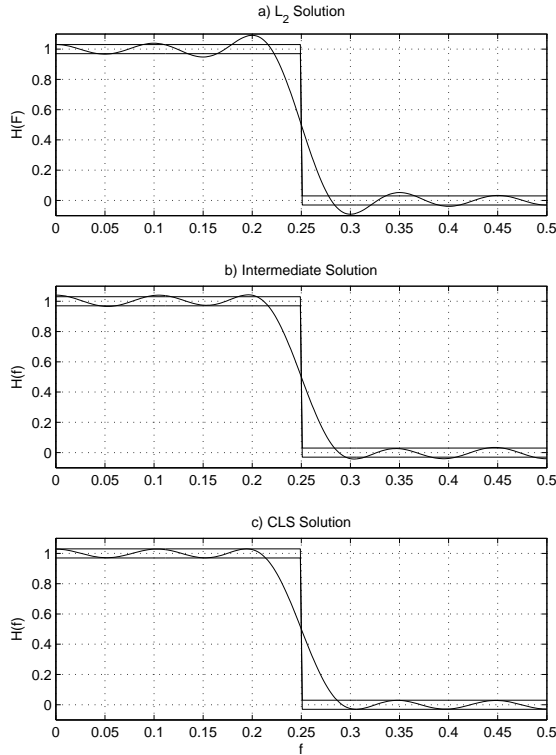


Fig. 14. CLS design example without transition bands.

design examples using the proposed approach. Figure 12 shows an example of a mild constraint ($\tau = 0.6$), whereas 13 illustrates an advantage of this method, associated to a hard constraint ($\tau = 0.3$). The method is trying iteratively to reduce the maximum error towards the constraint; however the specified constraint in Figure 13 is such that even at the point where an equiripple response is reached for the specified transition bands the constraint is not met. At this point the method converges to an optimal l_p solution that approximates equiripple as p increases (the examples provided use $p = 50$).

A different behavior occurs when no transition bands are defined. Departing from an initial l_2 guess (as shown in Figure 14.a) the proposed IRLS-based CLS algorithm begins weighting frequencies selectively in order to reduce the l_∞ error towards the constraints τ at each iteration. Eventually an equiripple behavior can be observed if the constraints are too harsh (as in Figure 14.b). The algorithm will keep weighting until all frequencies meet the constraints (as in Figure 14.c). The absence of a specified transition band presents some ambiguity in defining valid frequencies for weighting. One cannot (or rather should not) apply weights too close to the transition frequency specified as this would result in an effort by the algorithm to create a steep transition region (which as mentioned previously is counterintuitive to finding an equiripple solution). In a sense, this would mean having two opposite effects working at the same time and the algorithm cannot accommodate both, usually leading to numerical problems.

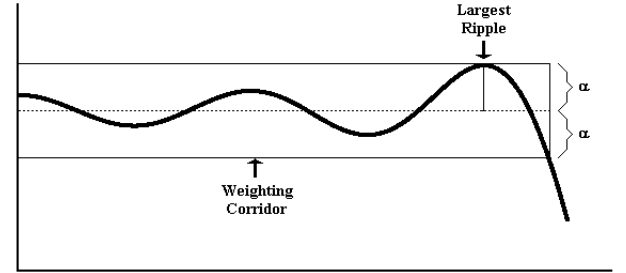


Fig. 15. Definition of *induced* transition band.

In order to avoid these issues, an algorithm can be devised that selects a subset of the sampled frequencies for weighting purposes at each iteration. The idea is to identify the largest ripple per band at each iteration (the ripple associated with the largest error for a given band) and select the frequencies within that band with errors equal or smaller than such ripple error. In this way one avoids weighting frequencies around the transition frequency. This idea is illustrated in Figure 15.

The previous example is fundamental since it illustrates the relevance of this method: since for a particular transition band the tightest constraint that one can get is given by the equiripple (or minimax) design (as shown in Section II-F1), a problem might arise when specifications are tighter than what the minimax design can meet. Adams found this problem (as reported in [16]); his method breaks under these conditions. The method proposed here overcomes an inadequate constraint and relaxes the transition band to meet the constraint.

It is worth noting that the polynomial weighting form works even when no transition bands are specified (this must become evident from Figure 14.c above). However, the user must be aware of some practical issues related to this approach. Figure 16 shows a typical CLS polynomial weighting function. Its "spiky" character becomes more dramatic as p increases (the method still follows the homotopy and partial updating ideas from previous sections) as shown in Figure

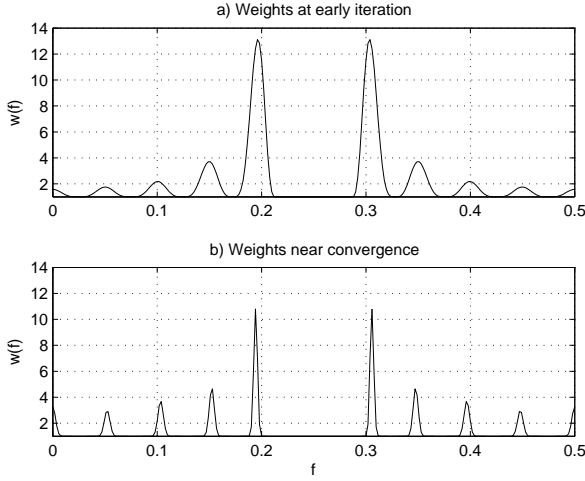


Fig. 16. CLS weights.

16.b. It must be evident that the algorithm will assign heavy weights to frequencies with large errors, but as p increases the difference in weighting exaggerates. At some point the user must make sure that proper sampling is done to ensure that frequencies with large weights (from a theoretical perspective) are being included in the problem, without compromising computational efficiency (by means of massive oversampling, which can lead to ill-conditioning in numerical least squares methods). Also as p increases, the range of frequencies with significantly large weights becomes narrower, thus reducing the overall weighting effect and affecting convergence speed.

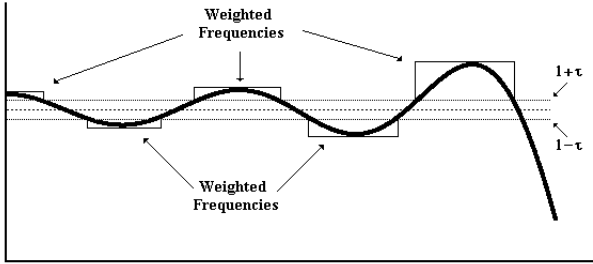


Fig. 17. CLS envelope weighting function.

A second weighting form can be defined where envelopes are used. The *envelope weighting function* approach works by assigning a weight to all frequencies not meeting a constraint. The value of such weights are assigned as flat intervals as illustrated in Figure 17. Intervals are determined by the edge frequencies within neighborhoods around peak error frequencies for which constraints are not met. Clearly these neighborhoods could change at each iteration. The weight of the k -th interval is still determined by our typical expression,

$$w_k(\omega) = |\epsilon(\omega_k^+)|^{\frac{p-2}{2}}$$

where ω_k^+ is the frequency with largest error within the k -th interval.

Envelope weighting has been applied in practice with good results. It is particularly effective at reaching high values of p without ill-conditioning, allowing for a true alternative to minimax design. Figure 18 shows an example using $\tau = 0.4$; the algorithm managed to find a solution for $p = 500$. By specifying transition bands and unachievable constraints one can produce an almost equiripple

solution in an efficient manner, with the added flexibility that milder constraints will result in CLS designs.

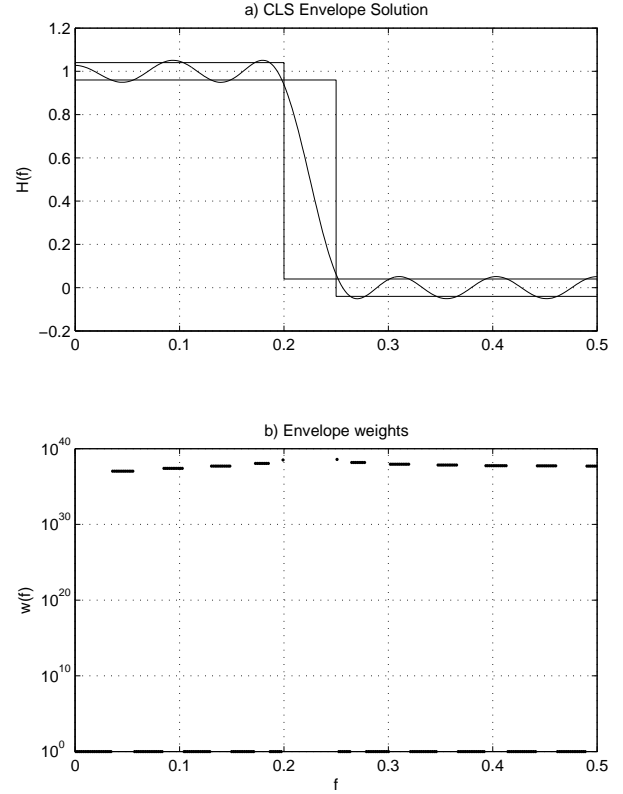


Fig. 18. CLS design example using envelope weights.

3) *Comparison with l_p problem:* This chapter presented two problems with similar effects. On one hand, Section II-B3 illustrated the fact (see Figure 6) that as p increases towards infinity, an l_p filter will approximate an l_∞ one. On the other hand, Section II-F presented the constrained least squared problem, and introduced IRLS-based algorithms that produce filters that approximate equiripple behavior as the constraint specifications tighten.

A natural question arises: how do these methods compare with each other? In principle it should be possible to compare their performances, as long as the necessary assumptions about the problem to be solved are compatible in both methods. Figure 19 shows a comparison of these algorithms with the following specifications:

- Both methods designed length-21 Type-I lowpass linear phase digital filters with fixed transition bands defined by $f = \{0.2, 0.24\}$ (in normalized linear frequency).
- The l_p experiment used the following values of $p = \{2, 2.2, 2.5, 3.4, 5, 7, 10, 15, 20, 30, 50, 70, 100, 170, 400\}$
- The CLS experiment used the polynomial weighting method with fixed transition bands and a value of $p = 60$. The error tolerances were $\tau = \{.06, .077, .078, .8, .084, .088, .093, .1, .11, .12, .13, .14, .15, .16, .17, .18\}$

Some conclusions can be derived from Figure 19. Even though at the extremes of the curves they both seem to meet, the CLS curve lies just below the l_p curve for most values of p and τ . These two facts should be expected: on one hand, in principle the CLS algorithm gives an l_2 filter if the constraints are so mild that they are not active for any frequency after the first iteration (hence the two curves should match

around $p = 2$). On the other hand, once the constraints become too harsh, the fixed transition band CLS method basically should design an equiripple filter, as only the active constraint frequencies are l_p -weighted (this effects is more noticeable with higher values of p). Therefore for tight constraints the CLS filter should approximate an l_∞ filter.

The reason why the CLS curve lies under the l_p curve is because for a given error tolerance (which could be interpreted as *for a given minimax error* ε_∞) the CLS method finds the optimal l_2 filter. An l_p filter is optimal in an l_p sense; it is not meant to be optimal in either the l_2 or l_∞ senses. Hence for a given τ it cannot beat the CLS filter in an l_2 sense (it can only match it, which happens around $p = 2$ or $p = \infty$).

It is important to note that both curves are not drastically different. While the CLS curve represents optimality in an $L_2 - l_\infty$ sense, not all the problems mentioned in this work can be solved using CLS filters (for example, the *magnitude* IIR problem presented in Section III-C2). Also, one of the objectives of this work is to motivate the use of l_p norms for filter design problems, and the proposed CLS implementations (which absolutely depends on IRLS-based l_p formulations) are good examples of the flexibility and value of l_p IRLS methods discussed in this work.

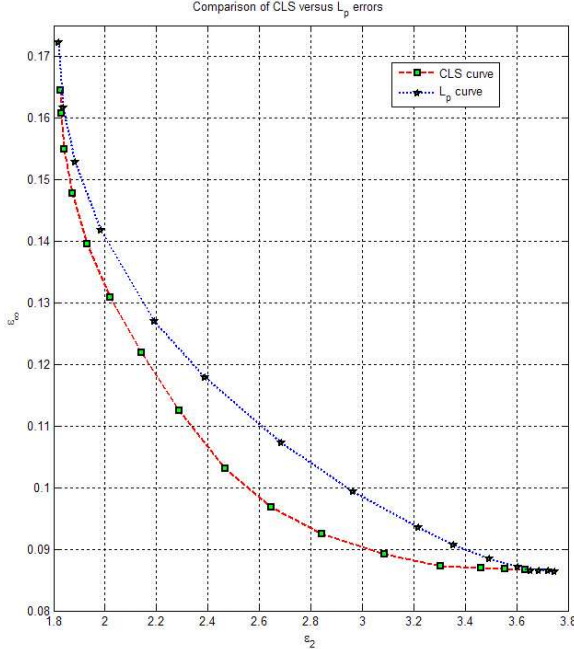


Fig. 19. Comparison between CLS and l_p problems.

III. INFINITE IMPULSE RESPONSE FILTERS

Chapter II introduced the problem of designing l_p FIR filters, along with several design scenarios and their corresponding design algorithms. This chapter considers the design of l_p IIR filters and examines the similarities and differences compared to l_p FIR filter design. It was mentioned in Section I-D that l_p FIR design involves a polynomial approximation. The problem becomes more complicated in the case of IIR filters as the approximation problem is a ratio of two polynomials. In fact, the case of FIR polynomial approximation is a special form of IIR rational approximation where the denominator is equal to 1.

Infinite Impulse Response (or *recursive*) digital filters constitute an important analysis tool in many areas of science (such as signal

processing, statistics and biology). The problem of designing IIR filters has been the object of extensive study. Several approaches are typically used in designing IIR filters, but a general procedure follows: given a desired filter specification (which may consist of an impulse response or a frequency specification), a predetermined approximation error criterion is optimized. Although one of the most widely used error criteria in Finite Impulse Response (FIR) filters is the *least-squares* criterion (which in most scenarios merely requires the solution of a linear system), least-squares (l_2) approximation for IIR filters requires an optimization over an infinite number of filter coefficients (in the time domain approximation case). Furthermore, optimizing for an IIR frequency response leads to a rational (nonlinear) approximation problem rather than the polynomial problem of FIR design.

As discussed in the previous chapter, a successful IRLS-based l_p algorithm depends to a large extent in the solution of a weighted l_2 problem. One could argue that one of the most important aspects contrasting FIR and IIR l_p filter design lies in the l_2 optimization step. This chapter presents the theoretical and computational issues involved in the design of both l_2 and l_p IIR filters and explores several approaches taken to handle the resulting nonlinear l_2 optimization problem. Section III-A introduces the IIR filter formulation and the nonlinear least-squares design problem. Section III-B presents the l_2 problem more formally, covering relevant methods as a manner of background and to lay down a framework for the approach proposed in this work. Some of the methods covered here date back to the 1960's, yet others are the result of current active work by a number of research groups; the approach employed in this work is described in section III-B13. Finally, Section III-C considers different design problems concerning IIR filters in an l_p sense, including IIR versions of the complex, frequency-varying and magnitude filter design problems as well as the proposed algorithms and their corresponding results.

A. IIR filters

An IIR filter describes a system with input $x(n)$ and output $y(n)$, related by the following expression

$$y(n) = \sum_{k=0}^M b(k)x(n-k) - \sum_{k=1}^N a(k)y(n-k)$$

Since the current output $y(n)$ depends on the input as well as on N previous output values, the output of an IIR filter might not be zero well after $x(n)$ becomes zero (hence the name "Infinite"). Typically IIR filters are described by a rational transfer function of the form

$$H(z) = \frac{B(z)}{A(z)} = \frac{b_0 + b_1 z^{-1} + \dots + b_M z^{-M}}{1 + a_1 z^{-1} + \dots + a_N z^{-N}} \quad (49)$$

where

$$H(z) = \sum_{n=0}^{\infty} h(n)z^{-n} \quad (50)$$

and $h(n)$ is the *infinite impulse response* of the filter. Its *frequency response* is given by

$$H(\omega) = H(z)|_{z=e^{j\omega}} \quad (51)$$

Substituting (49) into (51) we obtain

$$H(\omega) = \frac{B(\omega)}{A(\omega)} = \frac{\sum_{n=0}^M b_n e^{-j\omega n}}{1 + \sum_{n=1}^N a_n e^{-j\omega n}} \quad (52)$$

Given a *desired* frequency response $D(\omega)$, the l_2 IIR design problem consists of solving the following problem

$$\min_{a_n, b_n} \left| \frac{B(\omega)}{A(\omega)} - D(\omega) \right|_2^2 \quad (53)$$

for the $M + N + 1$ real filter coefficients a_n, b_n with $\omega \in \Omega$ (where Ω is the set of frequencies for which the approximation is done). A discrete version of (53) is given by

$$\min_{a_n, b_n} \sum_{\omega_k} \left| \frac{\sum_{n=0}^M b_n e^{-j\omega_k n}}{1 + \sum_{n=1}^N a_n e^{-j\omega_k n}} - D(\omega_k) \right|^2 \quad (54)$$

where ω_k are the L frequency samples over which the approximation is made. Clearly, (54) is a nonlinear least squares optimization problem with respect to the filter coefficients.

B. Least squares design of IIR filters

Section III-A introduced the IIR least squares design problem, as illustrated in (54). Such problem cannot be solved in the same manner as in the FIR case; therefore more sophisticated methods must be employed. As will be discussed later in Section III-C, some tradeoffs are desirable for l_p optimization. As in the case of FIR design, when designing l_p IIR filters one must use l_2 methods as internal steps over which one iterates while moving between different values of p . Clearly this internal iteration must not be too demanding computationally since an outer l_p loop will invoke it repeatedly (this process will be further illustrated in Section III-C1). With this issue in mind, one needs to select an l_2 algorithm that remains accurate within reasonable error bounds while remaining computationally efficient.

This section begins by summarizing some of the traditional approaches that have been employed for l_2 rational approximation, both within and outside filter design applications. Amongst the several existing traditional nonlinear optimization approaches, the Davidon-Fletcher-Powell (DFP) and the Gauss-Newton methods have been often used and remain relatively well understood in the filter design community. A brief introduction to both methods is presented in Section III-B1, and their caveats briefly explored.

An alternative to attacking a complex nonlinear problem like (54) with general nonlinear optimization tools consists in *linearization*, an attempt to "linearize" a nonlinear problem and to solve it by using linear optimization tools. Multiple efforts have been applied to similar problems in different areas of statistics and systems analysis and design. Section III-B2 introduces the notion of an *Equation Error*, a linear expression related to the actual *Solution Error* that one is interested in minimizing in l_2 design. The equation error formulation is nonetheless important for a number of filter design methods (including the ones presented in this work) such as Levy's method, one of the earliest and most relevant frequency domain linearization approaches. Section III-B4 presents a frequency domain equation error algorithm based on the methods by Prony and Padé. This algorithm illustrates the usefulness of the equation error formulation as it is fundamental to the implementation of the methods proposed later in this work (in Section III-C).

An important class of linearization methods fall under the name of *iterative prefiltering* algorithms, presented in Section III-B8. The *Sanathanan-Koerner (SK)* algorithm and the *Steiglitz-McBride (SMB)* methods are well known and commonly used examples in this category, and their strengths and weaknesses are explored. Another recent development in this area is the method by Jackson, also presented in this section. Finally, Soewito's *quasilinearization* (the method of choice for least squares IIR approximation in this work) is presented in Section III-B13.

1) *Traditional optimization methods*: One way to address (54) is to attempt to solve it with general nonlinear optimization tools. One of the most typical approach in nonlinear optimization is to apply either Newton's method or a Newton-based algorithm. One assumption of Newton's method is that the optimization function resembles a quadratic function near the solution. In order to update a current estimate, Newton's method requires first and second order information through the use of gradient and Hessian matrices. A *quasi-Newton method* is one that estimates in a certain way the second order information based on gradients (by generalizing the secant method to multiple dimensions).

One of the most commonly used quasi-Newton methods in IIR filter design is the *Davidon-Fletcher-Powell (DFP)* method [20]. In 1970 K. Steiglitz [49] used the DFP method to solve an IIR magnitude approximation to a desired real frequency response. For stability concerns he used a cascade form of the IIR filter given in (49) through

$$H(z) = \alpha \prod_{r=1}^M \frac{1 + a_r z^{-1} + b_r z^{-2}}{1 + c_r z^{-1} + d_r z^{-2}} \quad (55)$$

Therefore he considered the following problem,

$$\min_{a_n, b_n, c_n, d_n} \sum_{\omega_k} \left(\left| \alpha \prod_{r=1}^M \frac{1 + a_r e^{-j\omega_k} + b_r e^{-2j\omega_k}}{1 + c_r e^{-j\omega_k} + d_r e^{-2j\omega_k}} - D(\omega_k) \right|^2 \right)$$

His method is a direct implementation of the DFP algorithm in the problem described above.

In 1972 Andrew Deczky [50] employed the DFP algorithm to solve a complex IIR least- p approximation to a desired frequency response. Like Steiglitz, Deczky chose to employ the cascaded IIR structure of (55), mainly for stability reasons but also because he claims that for this structure it is simpler to derive the first order information required for the DFP method.

The MATLAB Signal Processing Toolbox includes a function called `INVREQZ`, originally written by J. Smith and J. Little [22]. `Invfreqz` uses the algorithm by Levy (see §III-B3) as an initial step and then begins an iterative algorithm based on the damped Gauss-Newton [21] to minimize the solution error ε_s according to the least-squared error criteria. This method performs a line search after every iteration to find the optimal direction for the next step. `Invfreqz` evaluates the roots of $A(z)$ after each iteration to verify that the poles of $H(z)$ lie inside the unit circle; otherwise it will convert the pole into its reciprocal. This approach guarantees a stable filter.

Among other Newton-based approaches, Spanos and Mingori [23] use a Newton algorithm combined with the Levenberg-Marquardt technique to improve the algorithm's convergence properties. Their idea is to express the denominator function $A(\omega)$ as a sum of second-order rational polynomials. Thus $H(\omega)$ can be written as

$$H(\omega) = \sum_{r=1}^{L-1} \frac{b_r + j\omega\beta_r}{a_r + j\omega\beta_r - \omega^2} + d$$

Their global descent approach is similar to the one presented in [24]. As any Newton-based method, this approach suffers under a poor initial guess, and does not guarantee to converge (if convergence occurs) to a local minimum. However, in such case, convergence is quadratic.

Kumaresan's method [25] considers a three-step approach. It is not clear whether his method attempts to minimize the equation error or the solution error. He uses divided differences [26] to reformulate the solution error in terms of the coefficients a_k . Using Lagrange multiplier theory, he defines

$$\mathcal{E} = \mathbf{y}^T \mathbf{C}^T [\mathbf{C} \mathbf{C}^T]^{-1} \mathbf{C} \mathbf{y} \quad (56)$$

where $\mathbf{y} = [H_0 H_1 \cdots H_{L-1}]^T$ contains the frequency samples and \mathbf{C} is a composition matrix containing the frequency divided differences and the coefficients a_k (a more detailed derivation can be found in [51]). Equation (56) is iterated until convergence of the coefficient vector $\hat{\mathbf{a}}$ is reached. This vector is used as initial guess in the second step, involving a Newton-Raphson search of the optimal $\hat{\mathbf{a}}$ that minimizes $\|\mathcal{E}\|_2$. Finally the vector $\hat{\mathbf{b}}$ is found by solving a linear system of equations.

2) *Equation error linearization methods*: Typically general use optimization tools prove effective in finding a solution. However in the context of IIR filter design, they often tend to take a rather large number of iterations, generate large matrices or require complicated steps like solving or estimating (and often inverting) vectors and matrices of first and second order information [35]. Using gradient-based tools for nonlinear problems like (54) certainly seems like a suboptimal approach. Also, typical Newton-based methods tend to converge quick (quadratically), yet they make assumptions about radii of convergence and initial proximity to the solution (otherwise performance is suboptimal). In the context of filter design one should wonder if better performance could be achieved by exploiting characteristics from the problem. This section introduces the concept of linearization, an alternative to general optimization methods that has proven successful in the context of rational approximation. The main idea behind linearization approaches consists in transforming a complex nonlinear problem into a sequence of linear ones, an idea that is parallel to the approach followed in our development of IRLS l_p optimization.

A common notion used in this work (as well as some of the literature related to linearization and filter design) is that there are two different error measures that authors often refer to. It is important to recognize the differences between them as one browses through literature. Typically one would be interested in minimizing the l_2 error given by:

$$\varepsilon = \|\mathcal{E}(\omega)\|_2^2 = \left\| D(\omega) - \frac{B(\omega)}{A(\omega)} \right\|_2^2 \quad (57)$$

This quantity is often referred to as the *solution error* (denoted by ε_s); we refer to the function $\mathcal{E}(\omega)$ in (57) as the *solution error function*, denoted by $\mathcal{E}_s(\omega)$. Also, in linearization algorithms the following measure often arises,

$$\varepsilon = \|\mathcal{E}(\omega)\|_2^2 = \|A(\omega)D(\omega) - B(\omega)\|_2^2 \quad (58)$$

This measure is often referred to as the *equation error* ε_e ; we denote the function $\mathcal{E}(\omega)$ in (58) as the *equation error function* $\mathcal{E}_e(\omega)$. Keeping the notation previously introduced, it can be seen that the two errors relate by one being a weighted version of the other,

$$\mathcal{E}_e(\omega) = A(\omega)\mathcal{E}_s(\omega)$$

3) *Levy's method*: E. C. Levy [27] considered in 1959 the following problem in the context of analog systems (electrical networks to be more precise): define¹

$$H(j\omega) = \frac{B_0 + B_1(j\omega) + B_2(j\omega)^2 + \cdots}{A_0 + A_1(j\omega) + A_2(j\omega)^2 + \cdots} = \frac{B(\omega)}{A(\omega)} \quad (59)$$

Given L samples of a desired complex-valued function $D(j\omega_k) = R(\omega_k) + jI(\omega_k)$ (where R, I are both real funtions of ω), Levy defines

$$\mathcal{E}(\omega) = D(j\omega) - H(j\omega) = D(j\omega) - \frac{B(\omega)}{A(\omega)}$$

or

$$\varepsilon = \sum_{k=0}^L |\mathcal{E}(\omega_k)|^2 = \sum_{k=0}^L |A(\omega_k)D(j\omega_k) - B(\omega_k)|^2 \quad (60)$$

Observing the linear structure (in the coefficients A_k, B_k) of equation (60), Levy proposed minimizing the quantity ε . He actually realized that this measure (what we would denote as the equation error) was indeed a **weighted** version of the actual solution error that one might be interested in; in fact, the denominator function $A(\omega)$ became the weighting function.

Levy's proposed method for minimizing (60) begins by writing ε as follows,

$$\varepsilon = \sum_{k=0}^L [(R_k \sigma_k - \omega_k \tau_k I_k - \alpha_k)^2 + (\omega_k \tau_k R_k + \sigma_k I_k - \omega_k \beta_k)^2] \quad (61)$$

by recognizing that (59) can be reformulated in terms of its real and imaginary parts,

$$H(j\omega) = \frac{\alpha + j\omega\beta}{\sigma + j\omega\tau}$$

with

$$\begin{aligned} \alpha + j\omega\beta &= (B_0 - B_2\omega^2 + B_4\omega^4 \cdots) \\ &\quad + j\omega(B_1 - B_3\omega^2 + B_5\omega^4 \cdots) \\ \sigma + j\omega\tau &= (A_0 - A_2\omega^2 + A_4\omega^4 \cdots) \\ &\quad + j\omega(A_1 - A_3\omega^2 + A_5\omega^4 \cdots) \end{aligned}$$

and performing appropriate manipulations². Note that the optimal set of coefficients A_k, B_k must satisfy

$$\frac{\partial \varepsilon}{\partial A_0} = \frac{\partial \varepsilon}{\partial A_1} = \cdots = \frac{\partial \varepsilon}{\partial B_0} = \cdots = 0$$

The conditions introduced above generate a linear system in the filter coefficients. Levy derives the system

$$\mathbf{C}\mathbf{x} = \mathbf{y} \quad (62)$$

where $\mathbf{C} = \{\mathbf{C}_1 \quad \mathbf{C}_2\}$ with

$$\mathbf{C}_1 = \begin{Bmatrix} \lambda_0 & 0 & -\lambda_2 & 0 & \lambda_4 & \cdots \\ 0 & \lambda_2 & 0 & -\lambda_4 & 0 & \cdots \\ \lambda_2 & 0 & -\lambda_4 & 0 & \lambda_6 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ T_1 & -S_2 & -T_3 & S_4 & T_5 & \cdots \\ S_2 & T_3 & -S_4 & -T_5 & S_6 & \cdots \\ T_3 & -S_4 & -T_5 & S_6 & T_7 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{Bmatrix}$$

$$\mathbf{C}_2 = \begin{Bmatrix} T_1 & S_2 & -T_3 & -S_4 & T_5 & \cdots \\ -S_2 & T_3 & S_4 & -T_5 & -S_6 & \cdots \\ T_3 & S_4 & -T_5 & -S_6 & T_7 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ U_2 & 0 & -U_4 & 0 & U_6 & \cdots \\ 0 & U_4 & 0 & -U_6 & 0 & \cdots \\ U_4 & 0 & -U_6 & 0 & U_8 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{Bmatrix}$$

¹For consistency with the rest of this document, notation has been modified from the author's original paper whenever deemed necessary.

²For further details on the algebraic manipulations involved, the reader should refer to [27].

and

$$\mathbf{x} = \begin{Bmatrix} B_0 \\ B_1 \\ B_2 \\ \vdots \\ A_1 \\ A_2 \\ \vdots \end{Bmatrix} \quad \mathbf{y} = \begin{Bmatrix} S_0 \\ T_1 \\ S_2 \\ T_3 \\ \vdots \\ 0 \\ U_2 \\ 0 \\ U_4 \\ \vdots \end{Bmatrix} \quad (63)$$

with

$$\begin{aligned} \lambda_h &= \sum_{l=0}^{L-1} \omega_l^h \\ S_h &= \sum_{l=0}^{L-1} \omega_l^h R_l \\ T_h &= \sum_{l=0}^{L-1} \omega_l^h I_l \\ U_h &= \sum_{l=0}^{L-1} \omega_l^h (R_l^2 + I_l^2) \end{aligned}$$

Solving for the vector \mathbf{x} from (62) gives the desired coefficients (note the trivial assumption that $A_0 = 1$). It is important to remember that although Levy's algorithm leads to a linear system of equations in the coefficients, his approach is indeed an equation error method. Matlab's `invfreqz` function uses an adaptation of Levy's algorithm for its least-squares equation error solution.

4) *Prony-based equation error linearization*: A number of algorithms that consider the approximation of functions in a least-squared sense using rational functions relate to Prony's method. This section summarizes these methods especially in the context of filter design.

5) *Prony's method*: The first method considered in this section is due to Gaspard Riche Baron de Prony, a Lyonnais mathematician and physicist which, in 1795, proposed to model the expansion properties of different gases by sums of damped exponentials. His method [29] approximates a sampled function $f(n)$ (where $f(n) = 0$ for $n < 0$) with a sum of N exponentials,

$$f(n) = \sum_{k=1}^N c_k e^{s_k n} = \sum_{k=1}^N c_k \lambda_k^n \quad (64)$$

where $\lambda_k = e^{s_k}$. The objective is to determine the N parameters c_k and the N parameters s_k in (64) given $2N$ samples of $f(n)$.

It is possible to express (64) in matrix form as follows,

$$\begin{bmatrix} 1 & 1 & \cdots & 1 \\ \lambda_1 & \lambda_2 & \cdots & \lambda_N \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^{N-1} & \lambda_2^{N-1} & \cdots & \lambda_N^{N-1} \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_N \end{bmatrix} = \begin{bmatrix} f(0) \\ f(1) \\ \vdots \\ f(N-1) \end{bmatrix} \quad (65)$$

System (65) has a Vandermonde structure with N equations, but $2N$ unknowns (both c_k and λ_k are unknown) and thus it cannot be solved directly. Yet the major contribution of Prony's work is to recognize that $f(n)$ as given in (64) is indeed the solution of a homogeneous order- N Linear Constant Coefficient Difference Equation (LCCDE) [52, ch. 4] given by

$$\sum_{p=0}^N a_p f(m-p) = 0 \quad (66)$$

with $a_0 = 1$. Since $f(n)$ is known for $0 \leq n \leq 2N-1$, we can extend (66) into an $(N \times N)$ system of the form

$$\begin{bmatrix} f(N-1) & f(N-2) & \cdots & f(0) \\ f(N) & f(N-1) & \cdots & f(1) \\ \vdots & \vdots & \ddots & \vdots \\ f(2N-2) & f(2N-3) & \cdots & f(N-1) \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_N \end{bmatrix} = \hat{\mathbf{f}} \quad (67)$$

where

$$\hat{\mathbf{f}} = \begin{bmatrix} -f(N) \\ -f(N+1) \\ \vdots \\ -f(2N-1) \end{bmatrix}$$

which we can solve for the coefficients a_p . Such coefficients are then used in the *characteristic equation* [53, §2.3] of (66),

$$\lambda^N + a_1 \lambda^{N-1} + \cdots + a_{N-1} \lambda + a_N = 0 \quad (68)$$

The N roots λ_k of (67) are called the *characteristic roots* of (66). From the λ_k we can find the parameters s_k using $s_k = \ln \lambda_k$. Finally, it is now possible to solve (65) for the parameters c_k .

The method described above is an adequate representation of Prony's original method [29]. More detailed analysis is presented in [54]–[57] and [58, §11.4]. Prony's method is an adequate algorithm for interpolating $2N$ data samples with N exponentials. Yet it is not a filter design algorithm as it stands. Its connection with IIR filter design, however, exists and will be discussed in the following sections.

6) *Padé's method*: The work by Prony served as inspiration to Henry Padé, a French mathematician which in 1892 published a work [30] discussing the problem of rational approximation. His objective was to approximate a function that could be represented by a power series expansion using a rational function of two polynomials.

Assume that a function $f(x)$ can be represented with a power series expansion of the form

$$f(x) = \sum_{k=0}^{\infty} c_k x^k \quad (69)$$

Padé's idea was to approximate $f(x)$ using the function

$$\hat{f}(x) = \frac{B(x)}{A(x)} \quad (70)$$

where

$$B(x) = \sum_{k=0}^M b_k x^k$$

and

$$A(x) = 1 + \sum_{k=1}^N a_k x^k$$

The objective is to determine the coefficients a_k and b_k so that the first $M+N+1$ terms of the residual

$$r(x) = A(x)f(x) - B(x)$$

dissappear (i.e. the first $N+M$ derivatives of $f(x)$ and $\hat{f}(x)$ are equal [59]). That is, [60],

$$r(x) = A(x) \sum_{k=0}^{\infty} c_k x^k - B(x) = x^{M+N+1} \sum_{k=0}^{\infty} d_k x^k$$

To do this, consider $A(x)f(x) = B(x)$ [56]

$$(1 + a_1 x + \cdots + a_N x^N) \cdot (c_0 + c_1 x + \cdots + c_i x^i + \cdots) = b_0 + b_1 x + \cdots + b_M x^M$$

By equating the terms with same exponent up to order $M + N + 1$, we obtain two sets of equations,

$$\begin{cases} c_0 = b_0 \\ a_1 c_0 + c_1 = b_1 \\ a_2 c_0 + a_1 c_1 + c_2 = b_2 \\ a_3 c_0 + a_2 c_1 + a_1 c_2 + c_3 = b_3 \\ \vdots \\ a_N c_{M-N} + a_{N-1} c_{M-N+1} + \cdots + c_M = b_M \end{cases} \quad (71)$$

$$\begin{cases} a_N c_{M-N+1} + a_{N-1} c_{M-N+2} + \cdots + c_{M+1} = 0 \\ a_N c_{M-N+2} + a_{N-1} c_{M-N+3} + \cdots + c_{M+2} = 0 \\ \vdots \\ a_N c_M + a_{N-1} c_{M+1} + \cdots + c_{M+N} = 0 \end{cases} \quad (72)$$

Equation (72) represents an $N \times N$ system that can be solved for the coefficients a_k given $c(n)$ for $0 \leq n \leq N + M$. These values can then be used in (71) to solve for the coefficients b_k . The result is a system whose impulse response matches the first $N + M + 1$ values of $f(n)$.

7) *Prony-based filter design methods*: Both the original methods by Prony and Pade were meant to interpolate data from applications that have little in common with filter design. What is relevant to this work is their use of rational functions of polynomials as models for data, and the linearization process they both employ.

When designing FIR filters, a common approach is to take L samples of the desired frequency response $D(\omega)$ and calculate the inverse DFT of the samples. This design approach is known as *frequency sampling*. It has been shown [28] that by designing a length- L filter $h(n)$ via the frequency sampling method and symmetrically truncating $h(n)$ to N values ($N \ll L$) it is possible to obtain a least-squares optimal length- N filter $h_N(n)$. It is not possible however to extend completely this method to the IIR problem. This section presents an extension based on the methods by Prony and Pade, and illustrates the shortcomings of its application.

Consider the frequency response defined in (52). One can choose L **equally spaced** samples of $H(\omega)$ to obtain

$$H(\omega_k) = H_k = \frac{B_k}{A_k} \quad \text{for } k = 0, 1, \dots, L-1 \quad (73)$$

where A_k and B_k represent the length- L DFTs of the filter coefficients a_n and b_n respectively. The division in (73) is done point-by-point over the L values of A_k and B_k . The objective is to use the relationship in described in (73) to calculate a_n and b_n .

One can express (73) as $B_k = H_k A_k$. This operation represents the length- L circular convolution $b(n) = h(n) \circledast a(n)$ defined as follows [43, §8.7.5]

$$b(n) = h(n) \circledast a(n) = \sum_{m=0}^{L-1} h[((n-m))_L] a(m), \quad 0 \leq n \leq L-1 \quad (74)$$

where $h(n)$ is the length- L inverse DFT of H_k and the operator $((\cdot))_L$ represents modulo L . Let

$$\hat{\mathbf{a}} = \begin{bmatrix} 1 \\ a_1 \\ \vdots \\ a_N \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{and} \quad \hat{\mathbf{b}} = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_M \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (75)$$

Therefore (74) can be posed as a matrix operation [28, §7.4.1] of the form

$$\hat{\mathbf{H}} \hat{\mathbf{a}} = \hat{\mathbf{b}} \quad (76)$$

where $\hat{\mathbf{H}} = [\hat{\mathbf{H}}_1 \hat{\mathbf{H}}_2]$ with

$$\hat{\mathbf{H}}_1 = \begin{bmatrix} h_0 & h_{L-1} & \cdots & h_{L-N} \\ h_1 & h_0 & \cdots & h_{L-N+1} \\ \vdots & \vdots & \ddots & \vdots \\ h_M & h_{M-1} & \cdots & h_{((L-N+M))_L} \\ h_{M+1} & h_M & \cdots & h_{((L-N+M+1))_L} \\ \vdots & \vdots & \ddots & \vdots \\ h_{L-2} & h_{L-3} & \cdots & h_{L-N-2} \\ h_{L-1} & h_{L-2} & \cdots & h_{L-N-1} \end{bmatrix}$$

$$\hat{\mathbf{H}}_2 = \begin{bmatrix} h_{L-N-1} & \cdots & h_2 & h_1 \\ h_{L-N} & \cdots & h_3 & h_2 \\ \vdots & \ddots & \vdots & \vdots \\ h_{((L-N+M-1))_L} & \cdots & h_{M+2} & h_{M+1} \\ h_{((L-N+M))_L} & \cdots & h_{M+3} & h_{M+2} \\ \vdots & \ddots & \vdots & \vdots \\ h_{L-N-3} & \cdots & h_0 & h_{L-1} \\ h_{L-N-2} & \cdots & h_1 & h_0 \end{bmatrix}$$

Hence $\hat{\mathbf{H}}$ is an $L \times L$ matrix. From (75) it is clear that the $L - (N + 1)$ rightmost columns of $\hat{\mathbf{H}}$ can be discarded (since the last $L - (N + 1)$ values of $\hat{\mathbf{a}}$ in (75) are equal to 0). Therefore equation (76) can be rewritten as

$$\begin{bmatrix} h_0 & h_{L-1} & \cdots & h_{L-N} \\ h_1 & h_0 & \cdots & h_{L-N+1} \\ \vdots & \vdots & \ddots & \vdots \\ h_M & h_{M-1} & \cdots & h_{((L-N+M))_L} \\ h_{M+1} & h_M & \cdots & h_{((L-N+M+1))_L} \\ \vdots & \vdots & \ddots & \vdots \\ h_{L-2} & h_{L-3} & \cdots & h_{L-N-2} \\ h_{L-1} & h_{L-2} & \cdots & h_{L-N-1} \end{bmatrix} \begin{bmatrix} 1 \\ a_1 \\ \vdots \\ a_N \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_M \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (77)$$

or in matrix notation

$$\mathbf{H} \begin{bmatrix} 1 \\ \mathbf{a} \end{bmatrix} = \begin{bmatrix} \mathbf{b} \\ \mathbf{0} \end{bmatrix} \quad \text{or} \quad \mathbf{H} \tilde{\mathbf{a}} = \tilde{\mathbf{b}} \quad (78)$$

where \mathbf{a} and \mathbf{b} correspond to the length- N and $(M + 1)$ filter coefficient vectors respectively and \mathbf{H} contains the first $N + 1$ columns of $\hat{\mathbf{H}}$. It is possible to uncouple the calculation of \mathbf{a} and \mathbf{b} from (78) by breaking \mathbf{H} furthermore as follows,

$$\mathbf{H} = \begin{bmatrix} \begin{matrix} h_0 & h_{L-1} & \cdots & h_{L-N} \\ h_1 & h_0 & \cdots & h_{L-N+1} \\ \vdots & \vdots & \ddots & \vdots \\ h_M & h_{M-1} & \cdots & h_{((L-N+M))_L} \end{matrix} \\ \begin{matrix} h_{M+1} & h_M & \cdots & h_{((L-N+M+1))_L} \\ \vdots & \vdots & \ddots & \vdots \\ h_{L-2} & h_{L-3} & \cdots & h_{L-N-2} \\ h_{L-1} & h_{L-2} & \cdots & h_{L-N-1} \end{matrix} \end{bmatrix}$$

Therefore

$$\mathbf{H} = \begin{bmatrix} \mathbf{H}_1 \\ \mathbf{H}_2 \end{bmatrix} \quad \tilde{\mathbf{a}} = \begin{bmatrix} \mathbf{b} \\ \mathbf{0} \end{bmatrix} \quad (79)$$

with

$$\tilde{\mathbf{a}} = \begin{bmatrix} 1 \\ \mathbf{a} \end{bmatrix}$$

as defined in (78). This formulation allows to uncouple the calculations for \mathbf{a} and \mathbf{b} using two systems,

$$\begin{aligned} \mathbf{H}_1 \tilde{\mathbf{a}} &= \mathbf{b} \\ \mathbf{H}_2 \tilde{\mathbf{a}} &= \mathbf{0} \end{aligned}$$

Note that the last equation can be expressed as

$$\hat{\mathbf{H}}_2 \mathbf{a} = -\hat{\mathbf{h}}_2 \quad (80)$$

where $\mathbf{H}_2 = [\hat{\mathbf{h}}_2 \ \hat{\mathbf{H}}_2]$ (that is, $\hat{\mathbf{h}}_2$ and $\hat{\mathbf{H}}_2$ contain the first and second through N -th columns of $\hat{\mathbf{H}}_2$ respectively).

From (80) one can conclude that if $L = N + M + 1$ and if $\hat{\mathbf{H}}_2$ and \mathbf{H}_1 are nonsingular, then they can be inverted³ to solve for the filter coefficient vectors \mathbf{a} in (80) and solve for \mathbf{b} using $\mathbf{H}_1 \tilde{\mathbf{a}} = \mathbf{b}$.

The algorithm described above is an **interpolation** method rather than an **approximation** one. If $L > N + M + 1$ and $\hat{\mathbf{H}}_2$ is full column rank then (80) is an overdetermined linear system for which no exact solution exists; therefore an approximation must be found. From (73) we can define the *solution error* function $\mathcal{E}_s(\omega_k)$ as

$$\mathcal{E}_s(\omega_k) = \frac{B(\omega_k)}{A(\omega_k)} - H(\omega_k) \quad (81)$$

Using this notation, the design objective is to solve the nonlinear problem

$$\min_{\mathbf{a}, \mathbf{b}} \|\mathcal{E}_s(\omega_k)\|_2^2$$

Consider the system in equation (78). If \mathbf{H}_2 is overdetermined, one can define an approximation problem by introducing an error vector \mathbf{e} ,

$$\hat{\mathbf{b}} = \mathbf{H} \tilde{\mathbf{a}} - \mathbf{e} \quad (82)$$

where

$$\mathbf{e} = \begin{bmatrix} e_1 \\ e_2 \end{bmatrix}$$

Again, it is possible to uncouple (82) as follows,

$$\mathbf{b} = \mathbf{H}_1 \tilde{\mathbf{a}} - \mathbf{e}_1 \quad (83)$$

$$\mathbf{e}_2 = \hat{\mathbf{h}}_2 + \hat{\mathbf{H}}_2 \mathbf{a} \quad (84)$$

One can minimize the least-squared error norm $\|\mathbf{e}_2\|_2$ of the overdetermined system (84) by solving the normal equations [21]

$$\hat{\mathbf{H}}_2^T \hat{\mathbf{h}}_2 = -\hat{\mathbf{H}}_2^T \hat{\mathbf{H}}_2 \mathbf{a}$$

so that

$$\mathbf{a} = -[\hat{\mathbf{H}}_2^T \hat{\mathbf{H}}_2]^{-1} \hat{\mathbf{H}}_2^T \hat{\mathbf{h}}_2$$

and use this result in (83)

$$\mathbf{b} = \mathbf{H}_1 \tilde{\mathbf{a}} \quad (85)$$

Equation (83) represents the following time-domain operation,

$$\varepsilon(n) = b(n) - h(n) \mathcal{D} a(n), \quad 0 \leq n \leq M$$

(where \mathcal{D} denotes *circular convolution*) and can be interpreted in the frequency domain as follows,

$$\mathcal{E}_e(\omega_k) = B(\omega_k) - H(\omega_k)A(\omega_k) \quad (86)$$

Equation (86) is a weighted version of (81), as follows

$$\mathcal{E}_e(\omega_k) = A(\omega_k) \mathcal{E}_s(\omega_k)$$

³In practice one should not invert the matrices \mathbf{H}_1 and $\hat{\mathbf{H}}_2$ but use a more robust and efficient algorithm. See [61] for details.

Therefore the algorithm presented above will find the filter coefficient vectors \mathbf{a} and \mathbf{b} that minimize the *equation error* \mathcal{E}_e in (86) in the least-squares sense. Unfortunately, this error is not what one may want to optimize, since it is a weighted version of the *solution error* \mathcal{E}_s .

8) *Iterative prefiltering linearization methods*: Section III-B2 introduced the equation error formulation and several algorithms that minimize it. In a general sense however one is more interested in minimizing the solution error problem from (54). This section presents several algorithms that attempt to minimize the solution error formulation from (53) by *prefiltering* the desired response $D(\omega)$ in (58) with $A(\omega)$. Then a new set of coefficients $\{a_n, b_n\}$ are found with an equation error formulation and the prefiltering step is repeated, hence defining an iterative procedure.

9) *Sanathanan-Koerner (SK) method*: The method by Levy presented in Section III-B3 suggests a relatively easy-to-implement approach to the problem of rational approximation. While interesting in itself, the equation error ε_e does not really represent what in principle one would like to minimize. A natural extension to Levy's method is the one proposed [31] by C. K. Sanathanan and J. Koerner in 1963. The algorithm iteratively *prefilters* the equation error formulation of Levy with an estimate of $A(\omega)$. The SK method considers the solution error function \mathcal{E}_s defined by

$$\begin{aligned} \mathcal{E}_s(\omega) &= D(\omega) - \frac{B(\omega)}{A(\omega)} = \frac{1}{A(\omega)} [A(\omega)D(\omega) - B(\omega)] \\ &= \frac{1}{A(\omega)} \mathcal{E}_e(\omega) \end{aligned} \quad (87)$$

Then the solution error problem can be written as

$$\min_{\mathbf{a}_k, \mathbf{b}_k} \varepsilon_s \quad (88)$$

where

$$\begin{aligned} \varepsilon_s &= \sum_{k=0}^L |\mathcal{E}_s(\omega_k)|^2 \\ &= \sum_{k=0}^L \frac{1}{|A(\omega_k)|^2} |\mathcal{E}_e(\omega_k)|^2 \\ &= W(\omega) |\mathcal{E}_e(\omega_k)|^2 \end{aligned} \quad (89)$$

Note that given $A(\omega)$, one can obtain an estimate for $B(\omega)$ by minimizing \mathcal{E}_e as Levy did. This approach provides an estimate, though, because one would need to know the optimal value of $A(\omega)$ to truly optimize for $B(\omega)$. The idea behind this method is that by solving iteratively for $A(\omega)$ and $B(\omega)$ the algorithm would eventually converge to the solution of the desired solution error problem defined by (88). Since $A(\omega)$ is not known from the beginning, it must be initialized with a reasonable value (such as $A(\omega_k) = 1$).

To solve (88) Sanathanan and Koerner defined the same linear system from (62) with the same matrix and vector definitions. However the scalar terms used in the matrix and vectors reflect the presence of the weighting function $W(\omega)$ in ε_s as follows,

$$\begin{aligned} \lambda_h &= \sum_{l=0}^{L-1} \omega_l^h W(\omega_l) \\ S_h &= \sum_{l=0}^{L-1} \omega_l^h R_l W(\omega_l) \\ T_h &= \sum_{l=0}^{L-1} \omega_l^h I_l W(\omega_l) \\ U_h &= \sum_{l=0}^{L-1} \omega_l^h (R_l^2 + I_l^2) W(\omega_l) \end{aligned}$$

Then, given an initial definition of $A(\omega)$, at the p -th iteration one sets

$$W(\omega) = \frac{1}{|A_{p-1}(\omega_k)|^2} \quad (90)$$

and solves (62) using $\{\lambda, S, T, U\}$ as defined above until a convergence criterion is reached. Clearly, solving (88) using (89) is equivalent to solving a series of weighted least squares problems where the weighting function consists of the estimated values of $A(\omega)$ from the previous iteration. This method is similar to a time-domain method proposed by Steiglitz and McBride [34], presented later in this chapter.

10) Method of Sid-Ahmed, Chottera and Jullien: The methods by Levy and Sanathanan and Koerner did arise from an analog analysis problem formulation, and cannot therefore be used directly to design digital filters. However these two methods present important ideas that can be translated to the context of filter design. In 1978 M. Sid-Ahmed, A. Chottera and G. Jullien followed on these two important works and adapted [32] the matrix and vectors used by Levy to account for the design of IIR digital filters, given samples of a desired frequency response. Consider the frequency response $H(\omega)$ defined in (52). In parallel with Levy's development, the corresponding equation error can be written as

$$\varepsilon_e = \sum_{k=0}^L |F_k(\omega)|^2 \quad (91)$$

with

$$F_k(\omega) = \left\{ (R_k + jI_k) \left(1 + \sum_{c=1}^N a_c e^{-j\omega_k c} \right) - \left(\sum_{c=0}^M b_c e^{-j\omega_k c} \right) \right\}$$

One can follow a similar differentiation step as Levy by setting

$$\frac{\partial \varepsilon_e}{\partial a_1} = \frac{\partial \varepsilon_e}{\partial a_2} = \dots = \frac{\partial \varepsilon_e}{\partial b_0} = \dots = 0$$

with as defined in (91). Doing so results in a linear system of the form

$$\mathbf{C}\mathbf{x} = \mathbf{y}$$

where the vectors \mathbf{x} and \mathbf{y} are given by

$$\mathbf{x} = \begin{Bmatrix} b_0 \\ \vdots \\ b_M \\ a_1 \\ \vdots \\ a_N \end{Bmatrix} \quad \mathbf{y} = \begin{Bmatrix} \phi_0 - r_0 \\ \vdots \\ \phi_M - r_M \\ -\beta_1 \\ \vdots \\ -\beta_N \end{Bmatrix} \quad (92)$$

The matrix \mathbf{C} has a special structure given by

$$\mathbf{C} = \begin{bmatrix} \mathbf{\Psi} & \mathbf{\Phi} \\ \mathbf{\Phi}^T & \mathbf{\Upsilon} \end{bmatrix}$$

where $\mathbf{\Psi}$ and $\mathbf{\Upsilon}$ are symmetric Toeplitz matrices of order $M+1$ and N respectively, and their first row is given by

$$\begin{aligned} \mathbf{\Psi}_{1m} &= \eta_{m-1} & \text{for } m = 1, \dots, M+1 \\ \mathbf{\Upsilon}_{1m} &= \beta_{m-1} & \text{for } m = 1, \dots, N \end{aligned}$$

Matrix $\mathbf{\Phi}$ has order $M+1 \times N$ and has the property that elements on a given diagonal are identical (i.e. $\mathbf{\Phi}_{i,j} = \mathbf{\Phi}_{i+1,j+1}$). Its entries are given by

$$\begin{aligned} \mathbf{\Phi}_{1m} &= \phi_m + r_m & \text{for } m = 1, \dots, N \\ \mathbf{\Phi}_{m1} &= \phi_{m-2} - r_{m-2} & \text{for } m = 2, \dots, M+1 \end{aligned}$$

The parameters $\{\eta, \phi, r, \beta\}$ are given by

$$\begin{aligned} \eta_i &= \sum_{k=0}^L \cos i\omega_k & \text{for } 0 \leq i \leq M \\ \beta_i &= \sum_{k=0}^L |D(\omega_k)|^2 \cos i\omega_k & \text{for } 0 \leq i \leq N-1 \\ \phi_i &= \sum_{k=0}^L R_k \cos i\omega_k & \text{for } 0 \leq i \leq \max(N, M-1) \\ r_i &= \sum_{k=0}^L I_k \sin i\omega_k & \text{for } 0 \leq i \leq \max(N, M-1) \end{aligned}$$

The rest of the algorithm works the same way as Levy's. For a solution error approach, one must weight each of the parameters mentioned above with the factor from (90) as in the SK method.

There are two important details worth mentioning at this point: on one hand the methods discussed up to this point (Levy, SK and Sid-Ahmed et al.) do not put any limitation on the spacing of the frequency samples; one can sample as fine or as coarse as desired in the frequency domain. On the other hand there is no way to decouple the solution of both numerator and denominator vectors. In other words, from (63) and (92) one can see that the linear systems that solve for vector \mathbf{x} solve for all the variables in it. This is more of an issue for the iterative methods (SK & Sid-Ahmed), since at each iteration one solves for all the variables, but for the purposes of updating one needs only to keep the denominator variables (they get used in the weighting function); the numerator variables are never used within an iteration (in contrast to Burrus' Prony-based method presented in Section III-B4). This approach *decouples* the numerator and denominator computation into two separate linear systems. One only needs to compute the denominator variables until convergence is reached, and only then it becomes necessary to compute the numerator variables. Therefore most of the iterations solve a smaller linear system than the methods involved up to this point.

11) Steiglitz-McBride iterative algorithm: In 1965 K. Steiglitz and L. McBride presented an algorithm [33], [34] that has become quite popular in statistics and engineering applications. The *Steiglitz-McBride* method (**SMB**) considers the problem of deriving a transfer function for either an analog or digital system from their input and output data; in essence it is a time-domain method. Therefore it is mentioned in this work for completeness as it closely relates to the methods by Levy, SK and Sid-Ahmed, yet it is far better known and understood.

The derivation of the SMB method follows closely that of SK. In the Z-domain, the transfer function of a digital system is defined by

$$H(z) = \frac{B(z)}{A(z)} = \frac{b_0 + b_1 z^{-1} + \dots + b_N z^{-N}}{1 + a_1 z^{-1} + \dots + a_N z^{-N}}$$

Furthermore

$$Y(z) = H(z)X(z) = \frac{B(z)}{A(z)}X(z)$$

Steiglitz and McBride define the following problem,

$$\min \varepsilon_s = \sum_i \mathcal{E}_i(z)^2 = \frac{1}{2\pi j} \oint \left| X(z) \frac{B(z)}{A(z)} - D(z) \right|^2 \frac{dz}{z} \quad (93)$$

where $X(z) = \sum_j x_j z^{-j}$ and $D(z) = \sum_j d_j z^{-j}$ represent the z-transforms of the input and desired signals respectively. Equation (93) is the familiar nonlinear solution error function expressed in the Z-domain. Steiglitz and McBride realized the complexity of such function and proposed the iterative solution (93) using a simpler

problem defined by

$$\min \varepsilon_e = \sum_i \mathcal{E}_i(z)^2 = \frac{1}{2\pi j} \oint |X(z)B(z) - D(z)A(z)|^2 \frac{dz}{z} \quad (94)$$

This linearized error function is the familiar equation error in the Z-domain. Steiglitz and McBride proposed a two-mode iterative approach. The **SMB Mode 1** iteration is similar to the SK method, in that at the k -th iteration a *linearized* error criterion based on (94) is used,

$$\begin{aligned} \mathcal{E}_k(z) &= \frac{B_k(z)}{A_{k-1}(z)} X(z) - \frac{A_k(z)}{A_{k-1}(z)} D(z) \\ &= W_k(z) [B_k(z)X(z) - A_k(z)D(z)] \end{aligned} \quad (95)$$

where

$$W_k(z) = \frac{1}{A_{k-1}(z)}$$

Their derivation⁴ leads to the familiar linear system

$$\mathbf{C}\mathbf{x} = \mathbf{y}$$

with the following vector definitions

$$\mathbf{x} = \begin{Bmatrix} b_0 \\ \vdots \\ b_N \\ a_1 \\ \vdots \\ a_N \end{Bmatrix} \quad \mathbf{q}_j = \begin{Bmatrix} x_j \\ \vdots \\ x_{j-N+1} \\ d_{j-1} \\ \vdots \\ d_{j-N} \end{Bmatrix}$$

The vector \mathbf{q}_j is referred to as the *input-output* vector. Then

$$\begin{aligned} \mathbf{C} &= \sum_j \mathbf{q}_j \mathbf{q}_j^T \\ \mathbf{y} &= \sum_j d_j \mathbf{q}_j \end{aligned}$$

SMB Mode 2 is an attempt at reducing further the error once Mode 1 produces an estimate close enough to the actual solution. The idea behind Mode 2 is to consider the solution error defined by (93) and equate its partial derivatives with respect to the coefficients to zero. Steiglitz and McBride showed [33], [34] that this could be attained by defining a new vector

$$\mathbf{r}_j = \begin{Bmatrix} x_j \\ \vdots \\ x_{j-N+1} \\ y_{j-1} \\ \vdots \\ y_{j-N} \end{Bmatrix}$$

Then

$$\begin{aligned} \mathbf{C} &= \sum_j \mathbf{r}_j \mathbf{q}_j^T \\ \mathbf{y} &= \sum_j d_j \mathbf{r}_j \end{aligned}$$

The main difference between Mode 1 and Mode 2 is the fact that Mode 1 uses the desired values to compute its vectors and matrices, whereas Mode 2 uses the actual output values from the filter. The rationale behind this is that at the beginning, the output function $y(t)$ is not accurate, so the desired function provides better data for computations. On the other hand, Mode 1 does not really solve the

desired problem. Once Mode 1 is deemed to have reached the vicinity of the solution, one can use true partial derivatives to compute the gradient and find the actual solution; this is what Mode 2 does.

It has been claimed that under certain conditions the Steiglitz-McBride algorithm converges. However no guarantee of global convergence exists. A more thorough discussion of the Steiglitz-McBride algorithm and its relationships to other parameter estimation algorithms (such as the Iterative Quadratic Maximum Likelihood algorithm, or IQML) are found in [62]–[64].

12) Jackson's method: The following is a recent approach (from 2008) by Leland Jackson [36] based in the frequency domain. Consider vectors $\mathbf{a} \in \mathbb{R}^N$ and $\mathbf{b} \in \mathbb{R}^M$ such that

$$H(\omega) = \frac{B(\omega)}{A(\omega)}$$

where $H(\omega)$, $B(\omega)$, $A(\omega)$ are the Fourier transforms of \mathbf{h} , \mathbf{b} and \mathbf{a} respectively. For a discrete frequency set one can describe Fourier transform vectors $\mathbf{B} = \mathbf{W}_b \mathbf{b}$ and $\mathbf{A} = \mathbf{W}_a \mathbf{a}$ (where $\mathbf{W}_b, \mathbf{W}_a$ correspond to the discrete Fourier kernels for \mathbf{b}, \mathbf{a} respectively). Define

$$H_a(\omega_k) = \frac{1}{A(\omega_k)}$$

In vector notation, let $\mathbf{D}_a = \text{diag}(\mathbf{H}_a) = \text{diag}(1/\mathbf{A})$. Then

$$H(\omega) = \frac{B(\omega)}{A(\omega)} = H_a(\omega)B(\omega) \Rightarrow \mathbf{H} = \mathbf{D}_a \mathbf{B} \quad (96)$$

Let $H_d(\omega)$ be the desired complex frequency response. Define $\mathbf{D}_d = \text{diag}(\mathbf{H}_d)$. Then one wants to solve

$$\min \mathbf{E}^* \mathbf{E} = \|\mathbf{E}\|_2^2$$

where $\mathbf{E} = \mathbf{H} - \mathbf{H}_d$. From (96) one can write $\mathbf{H} = \mathbf{H}_d + \mathbf{E}$ as

$$\mathbf{H} = \mathbf{D}_a \mathbf{B} = \mathbf{D}_a \mathbf{W}_b \mathbf{b} \quad (97)$$

Therefore

$$\mathbf{H}_d = \mathbf{H} - \mathbf{E} = \mathbf{D}_a \mathbf{W}_b \mathbf{b} - \mathbf{E} \quad (98)$$

Solving (98) for \mathbf{b} one gets

$$\mathbf{b} = (\mathbf{D}_a \mathbf{W}_b) \backslash \mathbf{H}_d \quad (99)$$

Also,

$$\mathbf{H}_d = \mathbf{D}_d \hat{\mathbf{I}} = \mathbf{D}_d \mathbf{D}_a \mathbf{A} = \mathbf{D}_a \mathbf{D}_d \mathbf{A} = \mathbf{D}_a \mathbf{D}_d \mathbf{W}_a \mathbf{a}$$

where $\hat{\mathbf{I}}$ is a unit column vector. Therefore

$$\mathbf{H} - \mathbf{E} = \mathbf{H}_d = \mathbf{D}_a \mathbf{D}_d \mathbf{W}_a \mathbf{a}$$

From (98) we get

$$\mathbf{D}_a \mathbf{W}_b \mathbf{b} - \mathbf{E} = \mathbf{D}_a \mathbf{D}_d \mathbf{W}_a \mathbf{a}$$

or

$$\mathbf{D}_a \mathbf{D}_d \mathbf{W}_a \mathbf{a} + \mathbf{E} = \mathbf{D}_a \mathbf{W}_b \mathbf{b}$$

which in a least squares sense results in

$$\mathbf{a} = (\mathbf{D}_a \mathbf{D}_d \mathbf{W}_a) \backslash (\mathbf{D}_a \mathbf{W}_b \mathbf{b}) \quad (100)$$

From (99) one gets

$$\mathbf{a} = (\mathbf{D}_a \mathbf{D}_d \mathbf{W}_a) \backslash (\mathbf{D}_a \mathbf{W}_b [(\mathbf{D}_a \mathbf{W}_b) \backslash \mathbf{H}_d])$$

As a summary, at the i -th iteration one can write (98) and (100) as follows,

$$\begin{aligned} \mathbf{b}_i &= (\text{diag}(1/\mathbf{A}_{i-1}) \mathbf{W}_b) \backslash \mathbf{H}_d \\ \mathbf{a}_i &= (\text{diag}(1/\mathbf{A}_{i-1}) \text{diag}(\mathbf{H}_d) \mathbf{W}_a) \backslash (\text{diag}(1/\mathbf{A}_{i-1}) \mathbf{W}_b \mathbf{b}_i) \end{aligned}$$

⁴For more details the reader should refer to [33], [34].

13) *Soewito's quasilinearization method*: Consider the equation error residual function

$$\begin{aligned}
e(\omega_k) &= B(\omega_k) - D(\omega_k) \cdot A(\omega_k) \\
&= \sum_{n=0}^M b_n e^{-j\omega_k n} - D(\omega_k) \cdot \left(1 + \sum_{n=1}^N a_n e^{-j\omega_k n}\right) \\
&= b_0 + b_1 e^{-j\omega_k} + \dots + b_M e^{-j\omega_k M} \dots \\
&\quad - D_k - D_k a_1 e^{-j\omega_k} - \dots - D_k a_N e^{-j\omega_k N} \\
&= (b_0 + \dots + b_M e^{-j\omega_k M}) - \dots \\
&\quad D_k (a_1 e^{-j\omega_k} + \dots + a_N e^{-j\omega_k N}) - D_k
\end{aligned}$$

with $D_k = D(\omega_k)$. The last equation indicates that one can represent the equation error in matrix form as follows,

$$\mathbf{e} = \mathbf{F}\mathbf{h} - \mathbf{D}$$

where $\mathbf{F} = [\mathbf{F}_1 \quad \mathbf{F}_2]$ and

$$\begin{aligned}
\mathbf{F}_1 &= \begin{bmatrix} 1 & e^{-j\omega_0} & \dots & e^{-j\omega_0 M} \\ \vdots & \vdots & & \vdots \\ 1 & e^{-j\omega_{L-1}} & \dots & e^{-j\omega_{L-1} M} \end{bmatrix} \\
\mathbf{F}_2 &= \begin{bmatrix} -D_0 e^{-j\omega_0} & \dots & -D_0 e^{-j\omega_0 N} \\ \vdots & & \vdots \\ -D_{L-1} e^{-j\omega_{L-1}} & \dots & -D_{L-1} e^{-j\omega_{L-1} N} \end{bmatrix}
\end{aligned}$$

and

$$\mathbf{h} = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_M \\ a_1 \\ \vdots \\ a_N \end{bmatrix} \quad \text{and} \quad \mathbf{D} = \begin{bmatrix} D_0 \\ \vdots \\ D_{L-1} \end{bmatrix}$$

Consider now the *solution error* residual function

$$\begin{aligned}
s(\omega_k) &= H(\omega_k) - D(\omega_k) = \frac{B(\omega_k)}{A(\omega_k)} - D(\omega_k) \\
&= \frac{1}{A(\omega_k)} [B(\omega_k) - D(\omega_k) \cdot A(\omega_k)] \\
&= W(\omega_k) e(\omega_k)
\end{aligned}$$

Therefore one can write the solution error in matrix form as follows

$$\mathbf{s} = \mathbf{W}(\mathbf{F}\mathbf{h} - \mathbf{D}) \quad (101)$$

where \mathbf{W} is a diagonal matrix with $\frac{1}{A(\omega)}$ in its diagonal. From (101) the least-squared solution error $\varepsilon_s = \mathbf{s}^* \mathbf{s}$ can be minimized by

$$\mathbf{h} = (\mathbf{F}^* \mathbf{W}^2 \mathbf{F})^{-1} \mathbf{F}^* \mathbf{W}^2 \mathbf{D} \quad (102)$$

From (102) an iteration⁵ could be defined as follows

$$\mathbf{h}_{i+1} = (\mathbf{F}^* \mathbf{W}_i^2 \mathbf{F})^{-1} \mathbf{F}^* \mathbf{W}_i^2 \mathbf{D}$$

by setting the weights \mathbf{W} in (101) equal to $A_k(\omega)$, the Fourier transform of the current solution for \mathbf{a} .

A more formal approach to minimizing ε_s consists in using a gradient method (these approaches are often referred to as *Newton-like* methods). First one needs to compute the *Jacobian* matrix \mathbf{J} of \mathbf{s} ,

⁵Soewito refers to this expression as the Steiglitz-McBride Mode-1 in frequency domain.

where the pq -th term of \mathbf{J} is given by $\mathbf{J}_{pq} = \frac{\partial s_p}{\partial h_q}$ with \mathbf{s} as defined in (101). Note that the p -th element of \mathbf{s} is given by

$$s_p = H_p - D_p = \frac{B_p}{A_p} - D_p$$

For simplicity one can consider these reduced form expressions for the independent components of \mathbf{h} ,

$$\begin{aligned}
\frac{\partial s_p}{\partial b_q} &= \frac{1}{A_p} \frac{\partial}{\partial b_q} \sum_{n=0}^M b_n e^{-j\omega_p n} = W_p e^{-j\omega_p q} \\
\frac{\partial s_p}{\partial a_q} &= B_p \frac{\partial}{\partial a_q} \frac{1}{A_p} = \frac{-B_p}{A_p^2} \frac{\partial}{\partial a_q} \left(1 + \sum_{n=1}^N a_n e^{-j\omega_p n}\right) \\
&= \frac{-1}{A_p} \cdot \frac{B_p}{A_p} \cdot e^{-j\omega_p q} = -W_p H_p e^{-j\omega_p q}
\end{aligned}$$

Therefore one can express the Jacobian \mathbf{J} as follows,

$$\mathbf{J} = \mathbf{W}\mathbf{G} \quad (103)$$

where $\mathbf{G} = [\mathbf{G}_1 \quad \mathbf{G}_2]$ and

$$\begin{aligned}
\mathbf{G}_1 &= \begin{bmatrix} 1 & e^{-j\omega_0} & \dots & e^{-j\omega_0 M} \\ \vdots & \vdots & & \vdots \\ 1 & e^{-j\omega_{L-1}} & \dots & e^{-j\omega_{L-1} M} \end{bmatrix} \\
\mathbf{G}_2 &= \begin{bmatrix} -H_0 e^{-j\omega_0} & \dots & -H_0 e^{-j\omega_0 N} \\ \vdots & & \vdots \\ -H_{L-1} e^{-j\omega_{L-1}} & \dots & -H_{L-1} e^{-j\omega_{L-1} N} \end{bmatrix}
\end{aligned}$$

Consider the *solution error* least-squares problem given by

$$\min_{\mathbf{h}} f(\mathbf{h}) = \mathbf{s}^T \mathbf{s}$$

where \mathbf{s} is the solution error residual vector as defined in (101) and depends on \mathbf{h} . It can be shown [21, pp. 219] that the gradient of the squared error $f(\mathbf{h})$ (namely ∇f) is given by

$$\nabla f = \mathbf{J}^* \mathbf{s} \quad (104)$$

A necessary condition for a vector \mathbf{h} to be a local minimizer of $f(\mathbf{h})$ is that the gradient ∇f be zero at such vector. With this in mind and combining (101) and (103) in (104) one gets

$$\nabla f = \mathbf{G}^* \mathbf{W}^2 (\mathbf{F}\mathbf{h} - \mathbf{D}) = \mathbf{0} \quad (105)$$

Solving the system (105) gives

$$\mathbf{h} = (\mathbf{G}^* \mathbf{W}^2 \mathbf{F})^{-1} \mathbf{G}^* \mathbf{W}^2 \mathbf{D}$$

An iteration can be defined as follows⁶

$$\mathbf{h}_{i+1} = (\mathbf{G}_i^* \mathbf{W}_i^2 \mathbf{F})^{-1} \mathbf{G}_i^* \mathbf{W}_i^2 \mathbf{D} \quad (106)$$

where matrices \mathbf{W} and \mathbf{G} reflect their dependency on current values of \mathbf{a} and \mathbf{b} .

Atmadji Soewito [35] expanded the method of *quasilinearization* of Bellman and Kalaba [65] to the design of IIR filters. To understand his method consider the first order of Taylor's expansion near $H_i(z)$, given by

$$\begin{aligned}
H_{i+1}(z) &= H_i(z) + \frac{[B_{i+1}(z) - B_i(z)]A_i(z) - [A_{i+1}(z) - A_i(z)]B_i(z)}{A_i^2(z)} \\
&= H_i(z) + \frac{B_{i+1}(z) - B_i(z)}{A_i(z)} - \frac{B_i(z)[A_{i+1}(z) - A_i(z)]}{A_i^2(z)}
\end{aligned}$$

⁶Soewito refers to this expression as the Steiglitz-McBride Mode-2 in frequency domain. Compare to the Mode-1 expression and the use of G_i instead of F .

Using the last result in the solution error residual function $s(\omega)$ and applying simplification leads to

$$\begin{aligned} s(\omega) &= \frac{B_{i+1}(\omega)}{A_i(\omega)} - \frac{H_i(\omega)A_{i+1}(\omega)}{A_i(\omega)} + \frac{B_i(\omega)}{A_i(\omega)} - D(\omega) \\ &= \frac{1}{A_i(\omega)} [B_{i+1}(\omega) - H_i(\omega)A_{i+1}(\omega) + B_i(\omega) \dots \\ &\quad - A_i(\omega)D(\omega)] \end{aligned} \quad (107)$$

Equation (107) can be expressed (dropping the use of ω for simplicity) as

$$s = W \left(([B_{i+1} - H_i(A_{i+1} - 1)] - H_i) + ([B_i - D(A_i - 1)] - D) \right) \quad (108)$$

One can recognize the two terms in brackets as $\mathbf{G}h_{i+1}$ and $\mathbf{F}h_i$ respectively. Therefore (108) can be represented in matrix notation as follows,

$$\mathbf{s} = \mathbf{W}[\mathbf{G}h_{i+1} - (\mathbf{D} + \mathbf{H}_i - \mathbf{F}h_i)] \quad (109)$$

with $\mathbf{H} = [H_0, H_1, \dots, H_{L-1}]^T$. Therefore one can minimize $\mathbf{s}^T \mathbf{s}$ from (109) with

$$h_{i+1} = (\mathbf{G}_i^* \mathbf{W}_i^2 \mathbf{G}_i)^{-1} \mathbf{G}_i^* \mathbf{W}_i^2 (\mathbf{D} + \mathbf{H}_i - \mathbf{F}h_i) \quad (110)$$

since all the terms inside the parenthesis in (110) are constant at the $(i+1)$ -th iteration. In a sense, (110) is similar to (106), where the desired function is updated from iteration to iteration as in (110).

It is important to note that any of the three algorithms can be modified to solve a *weighted* l_2 IIR approximation using a weighting function $W(\omega)$ by defining

$$V(\omega) = \frac{W(\omega)}{A(\omega)} \quad (111)$$

Taking (111) into account, the following is a summary of the three different updates discussed so far:

$$\text{SMB Frequency Mode-1: } h_{i+1} = (\mathbf{F}^* \mathbf{V}_i^2 \mathbf{F})^{-1} \mathbf{F}^* \mathbf{V}_i^2 \mathbf{D}$$

$$\text{SMB Frequency Mode-2: } h_{i+1} = (\mathbf{G}_i^* \mathbf{V}_i^2 \mathbf{F})^{-1} \mathbf{G}_i^* \mathbf{V}_i^2 \mathbf{D}$$

$$\text{Soewito's quasilinearization: } h_{i+1} = (\mathbf{G}_i^* \mathbf{V}_i^2 \mathbf{G}_i)^{-1} * \dots \\ \mathbf{G}_i^* \mathbf{V}_i^2 (\mathbf{D} + \mathbf{H}_i - \mathbf{F}h_i)$$

C. l_p approximation

Infinite Impulse Response (IIR) filters are important tools in signal processing. The flexibility they offer with the use of poles and zeros allows for relatively small filters meeting specifications that would require somewhat larger FIR filters. Therefore designing IIR filters in an efficient and robust manner is an important problem.

This section covers the design of a number of important l_p IIR problems. The methods proposed are consistent with the methods presented for FIR filters, allowing one to build up on the lessons learned from FIR design problems. The complex l_p IIR problem is first presented in Section III-C1, being an essential tool for other relevant problems. The l_p frequency-dependent IIR problem is also introduced in Section III-C1. While the frequency-dependent formulation might not be practical in itself as a filter design formulation, it is fundamental for the more relevant magnitude l_p IIR filter design problem, presented in Section III-C2.

Some complications appear when designing IIR filters, among which the intrinsic least squares solving step clearly arises from the rest. Being a nonlinear problem, special handling of this step is required. It was determined after thorough experimentation that the *quasilinearization* method of Soewito presented in Section III-B13 can be employed successfully to handle this issue.

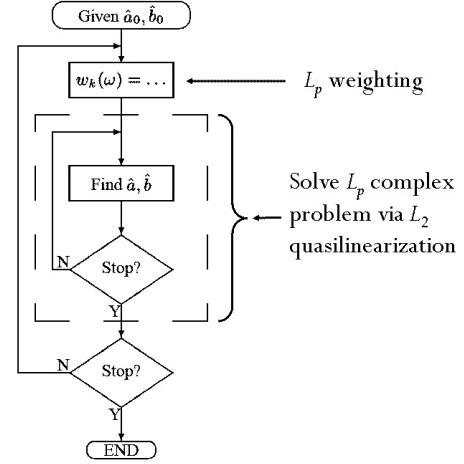


Fig. 20. Block diagram for complex l_p IIR algorithm.

1) *Complex and frequency-dependent l_p approximation:* Chapter II introduced the problem of designing l_p complex FIR filters. The complex l_p IIR algorithm builds up on its FIR counterpart by introducing a *nested structure* that internally solves for an l_2 complex IIR problem. Figure 20 illustrates this procedure in more detail. This method was first presented in [66].

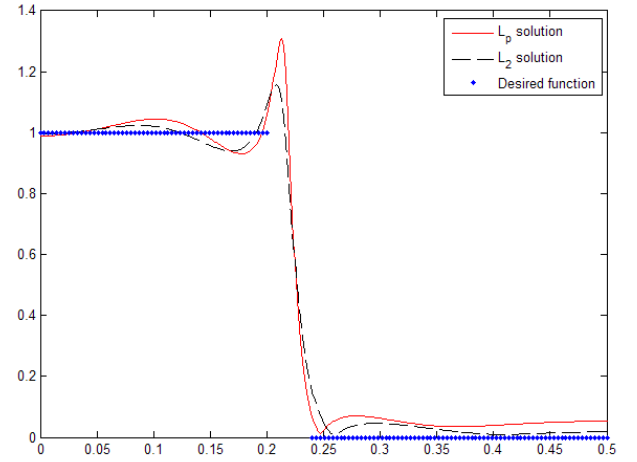


Fig. 21. Results for complex l_{100} IIR design.

Compared to its FIR counterpart, the IIR method only replaces the weighted linear least squares problem for Soewito's quasilinearization algorithm. While this nesting approach might suggest an increase in computational expense, it was found in practice that after the initial l_2 iteration, in general the l_p iterations only require from one to only a few internal weighted l_2 quasilinearization iterations, thus maintaining the algorithm efficiency. Figures 21 through 23 present results for a design example using a length-5 IIR filter with $p = 100$ and transition edge frequencies of 0.2 and 0.24 (in normalized frequency).

Figure 21 compares the l_2 and l_p results and includes the desired frequency samples. Note that no transition band was specified. Figure 22 illustrates the effect of increasing p . The largest error for the l_2 solution is located at the transition band edges. As p increases the algorithm weights the larger errors heavier; as a result the largest errors tend to decrease. In this case the magnitude of the frequency

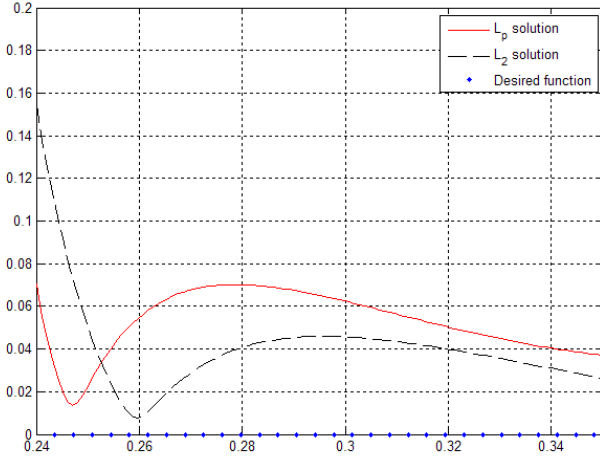


Fig. 22. Maximum error for l_2 and l_{100} complex IIR designs.

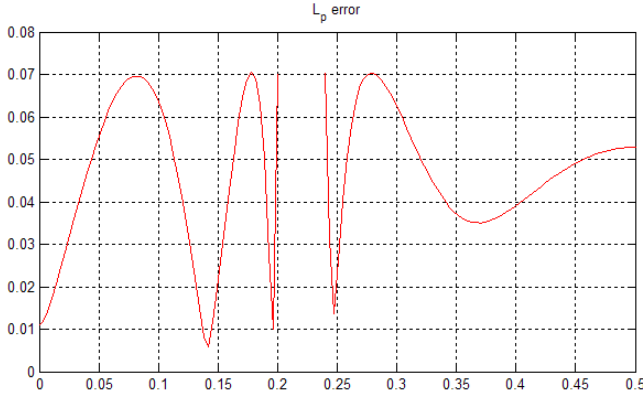


Fig. 23. Error curve for l_{100} complex IIR design.

response went from 0.155 at the stopband edge (in the l_2 case) to 0.07 (for the l_p design). Figure 23 shows the error function for the l_p design, illustrating the quasiequiripple behavior for large values of p .

Another fact worth noting from Figure 21 is the increase in the peak in the right hand side of the passband edge (around $f = 0.22$). The l_p solution increased the amplitude of this peak with respect to the corresponding l_2 solution. This is to be expected, since this peak occurs at frequencies not included in the specifications, and since the l_p algorithm will move poles and zeros around in order to meet find the optimal l_p solution (based on the frequencies included for the filter derivation). The addition of a specified transition band function (such as a spline) would allow for control of this effect, depending on the user's preferences.

The frequency-dependent FIR problem was first introduced in Section II-E. Following the FIR approach, one can design IIR frequency-dependent filters by merely replacing the linear weighted least squares step by a nonlinear approach, such as the quasilinearization method presented in Section III-B13 (as in the complex l_p IIR case). This problem illustrates the flexibility in design for l_p IRLS-based methods.

2) *Magnitude l_p IIR design*: The previous sections present algorithms that are based on complex specifications; that is, the user must specify both desired magnitude and phase responses. In some cases it might be better to specify a desired magnitude response only, while

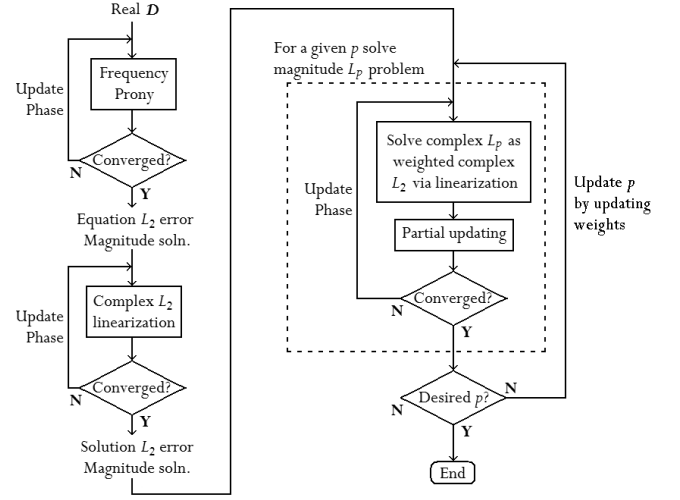


Fig. 24. Block diagram for magnitude l_p IIR method.

allowing an algorithm to select the phase that optimally minimizes the magnitude error. Note that if an algorithm is given a phase in addition to a magnitude function, it must then make a compromise between approximating both functions. The magnitude l_p IIR approximation problem overcomes this dilemma by posing the problem only in terms of a desired magnitude function. The algorithm would then find the optimal phase that provides the optimal magnitude approximation. A mathematical formulation follows,

$$\min_{\mathbf{a}, \mathbf{b}} \left\| |D(\omega)| - \left| \frac{B(\omega; \mathbf{b})}{A(\omega; \mathbf{a})} \right| \right\|_p^p \quad (112)$$

A critical idea behind the magnitude approach is to allow the algorithm to find the optimum phase for a magnitude approximation. It is important to recognize that the optimal magnitude filter indeed has a complex frequency response. Atmadji Soewito [35] published in 1990 a theorem in the context of l_2 IIR design that demonstrated that the phase corresponding to an optimal magnitude approximation could be found iteratively by **updating the desired phase** in a complex approximation scenario. In other words, given a desired complex response D_0 one can solve a complex l_2 problem and take the resulting phase to form a new desired response D^+ from the original desired magnitude response with the new phase. That is,

$$D_{i+1} = |D_0|e^{j\phi_i}$$

where D_0 represents the original desired magnitude response and $e^{j\phi_i}$ is the resulting phase from the previous iteration. This approach was independently suggested [36] by Leland Jackson and Stephen Kay in 2008.

This work introduces an algorithm to solve the magnitude l_p IIR problem by combining the IRLS-based complex l_p IIR algorithm from Section III-C1 with the phase updating ideas from Soewito, Jackson and Kay. The resulting algorithm is robust, efficient and flexible, allowing for different orders in the numerator and denominator as well as even or uneven sampling in frequency space, plus the optional use of specified transition bands. A block diagram for this method is presented in Figure 24.

The overall l_p IIR magnitude procedure can be summarized as follows,

- 1) Experimental analysis demonstrated that a reasonable initial solution for each of the three main stages would allow for faster convergence. It was found that the frequency domain Prony method by Burrus [28] (presented in Section III-B4) offered

a good *initial guess*. In Figure 24 this method is iterated to update the specified phase. The outcome of this step would be an **equation error l_2 magnitude** design.

- 2) The equation error l_2 magnitude solution from the previous step initializes a second stage where one uses quasilinearization to update the desired phase. Quasilinearization solves the true *solution error* complex approximation. Therefore by iterating on the phase one finds at convergence a **solution error l_2 magnitude** design.
- 3) The rest of the algorithm follows the same idea as in the previous step, except that the least squared step becomes a *weighted* one (to account for the necessary l_p homotopy weighting). It is also crucial to include the partial updating introduced in Section II-B2. By iterating on the weights one would find a **solution error l_p magnitude** design.

Figures 25 through 29 illustrate the effectiveness of this algorithm at each of the three different stages for length-5 filters **a** and **b**, with transition edge frequencies of 0.2 and 0.24 (in normalized frequency) and $p = 30$. A linear transition band was specified. Figures 25, 25 show the equation error l_2 , solution error l_2 and solution error l_p . Figure 28 shows a comparison of the magnitude error functions for the solution error l_2 and l_p designs. Figure 29 shows the phase responses for the three designs.

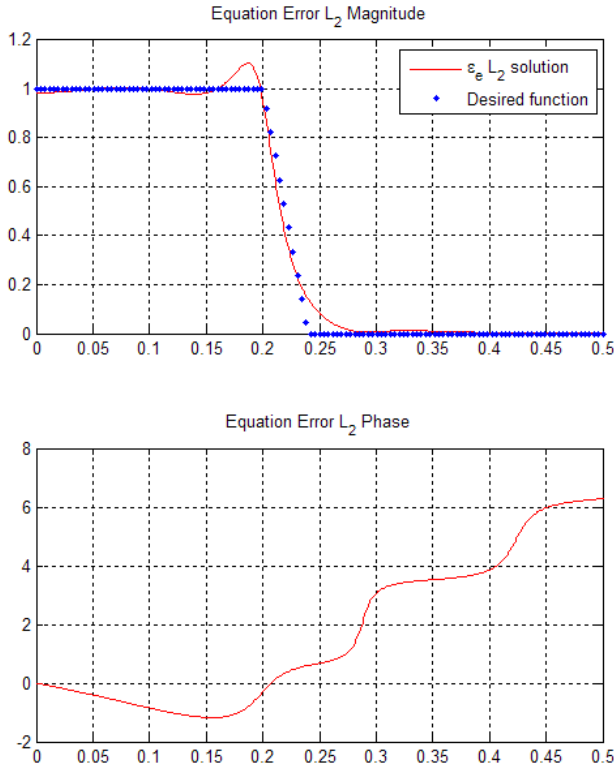


Fig. 25. Equation error l_2 magnitude design.

From Figures 28 and 29 one can see that the algorithm has changed the phase response in a way that makes the maximum magnitude error (located in the stopband edge frequency) to be reduced by approximately half its value. Furthermore, Figure 28 demonstrates that one can reach quasiequiripple behavior with relatively low values of p (for the examples shown, p was set to 30).

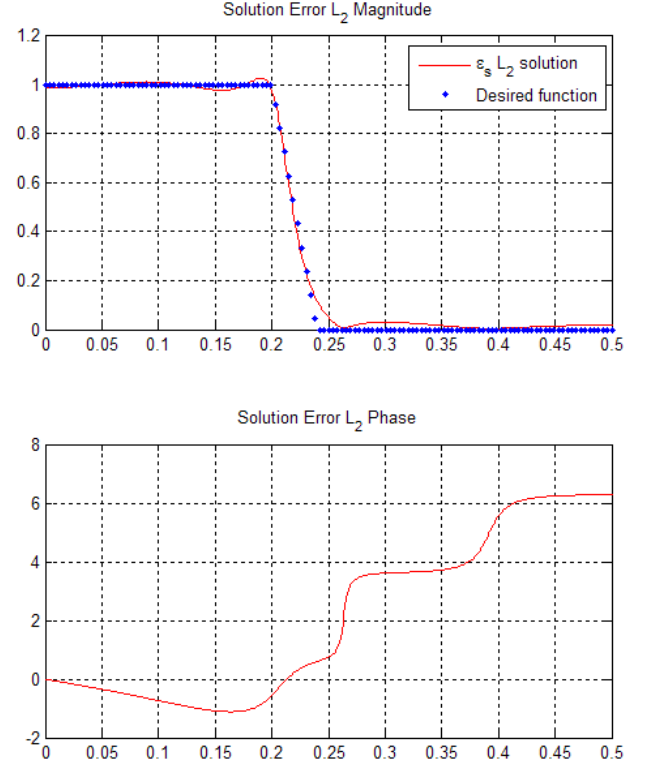


Fig. 26. Solution error l_2 magnitude design.

IV. CONCLUSIONS

Digital filters are essential building blocks for signal processing applications. One of the main goals of this work is to illustrate the versatility and relevance of l_p norms in the design of digital filters. While popular and well understood, l_2 and l_∞ filters do tend to accentuate specific benefits from their respective designs; filters designed using l_p norms as optimality criteria can offer a tradeoff between the benefits of these two commonly used criteria. This work presented a number of applications of L_p norms in both FIR and IIR filter design, and their corresponding design algorithms and software implementation.

The basic workhorse for the methods presented in this document is the *Iterative Reweighted Least Squares* algorithm, a simple yet powerful method that sets itself naturally adept for the design of l_p filters. The notion of converting a mathematically complex problem into a series of significantly easier optimization problems is common in optimization. Nevertheless, the existence results from Theorem 1 strongly motivate the use of IRLS methods to design l_p filters. Knowing that optimal weights exist that would turn the solution of a weighted least squared problem into the solution of a least- p problem must at the very least captivate the curiosity of the reader. The challenge lies in finding a robust and efficient method to find such weights. All the methods presented in this work work under this basic framework, updating iteratively the weighting function of a least squares problem in order to find the optimal l_p filter for a given application. Therefore it is possible to develop a suite of computer programs in a modular way, where with few adjustments one can solve a variety of problems.

Throughout this document one can find examples of the versatility of the IRLS approach. One can change the internal linear objective

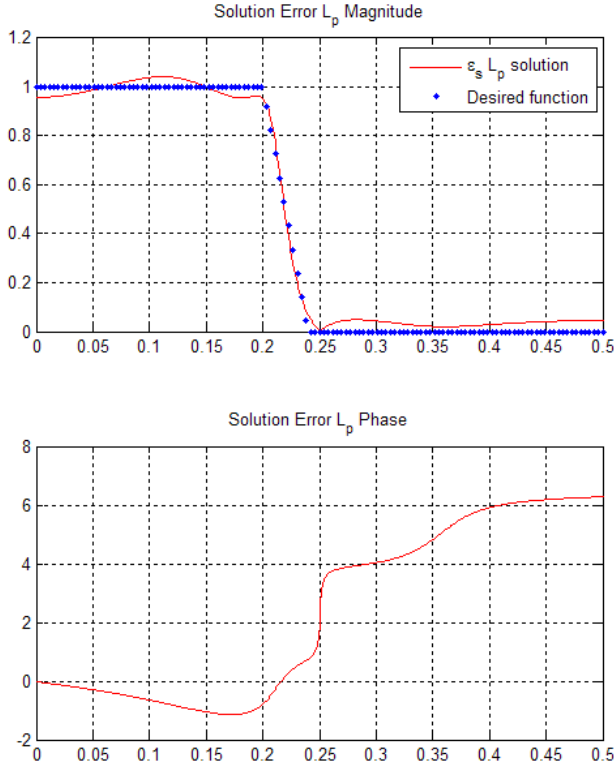


Fig. 27. Solution error l_p magnitude design.

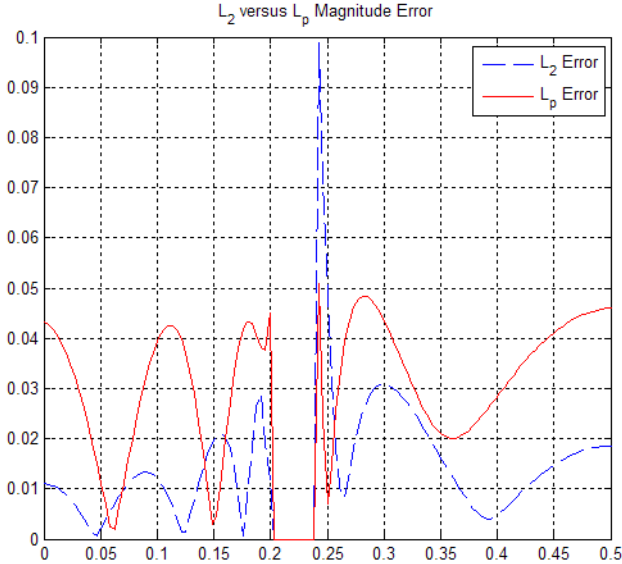


Fig. 28. Comparison of l_2 and l_p IIR magnitude designs

function from a complex exponential kernel to a sinusoidal one to solve complex and linear phase FIR filters respectively using the same algorithm. Further adaptations can be incorporated with ease, such as the proposed *adaptive solution* to improve robustness.

Another important design example permits to make p into a function of frequency to allow for different p -norms in different

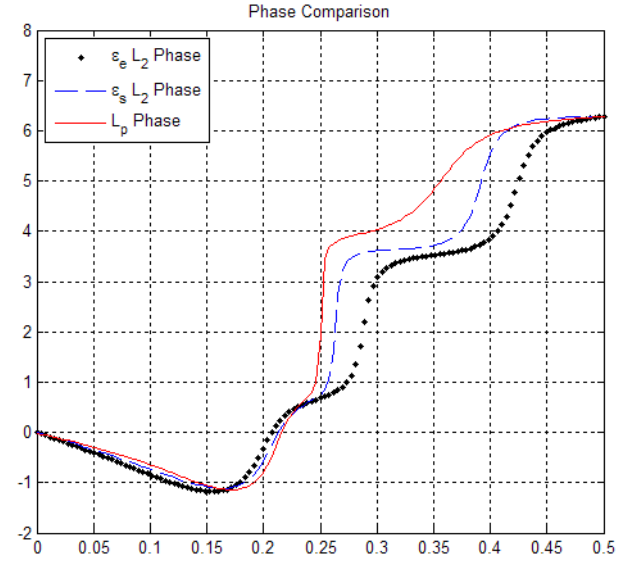


Fig. 29. Phase responses for l_2 and l_p IIR magnitude designs.

frequency bands. Such design merely requires a few changes in the implementation of the algorithm, yet allows for fancier, more elegant problems to be solved, such as the *Constrained Least Squares* (CLS) problem. In the context of FIR filters, this document presents the CLS problem from an l_p perspective. While the work by John Adams [16] set a milestone in digital filter design, this work introduces a strong algorithm and a different perspective to the problem from that by Adams and other authors. The IRLS l_p -based approach from this work proves to be robust and flexible, allowing for even and uneven sampling. Furthermore, while a user can use fixed transition bands, one would benefit much from using a flexible transition band formulation, where the proposed IRLS-based algorithm literally finds the optimal transition band definition based on the constraint specifications. Such flexibility allows for tight constraints that would otherwise cause other algorithms to fail meeting the constraint specifications, or simply not converging at all. Section II-F introduced two problem formulations as well as results that illustrate the method's effectiveness at solving the CLS problem.

While previous work exists in the area of FIR design (or in linear l_p approximation for that matter), the problem of designing l_p IIR filters has been far less explored. A natural reason for this is the fact that l_2 IIR design is in itself an open research area (and a rather complicated problem as well). Traditional linear optimization approaches cannot be directly used for either of these problems, and nonlinear optimization tools often prove either slow or do not converge.

This work presents the l_p IIR design problem as a natural extension of the FIR counterpart, where in a modular fashion the linear weighted l_2 section of the algorithms is replaced by a nonlinear weighted l_2 version. This problem formulation allows for the IIR implementation of virtually all the IRLS FIR methods presented in Chapter II. Dealing with the weighted nonlinear l_2 problem is a different story.

The problem of rational l_2 approximation has been studied for some time. However the sources of ideas and results related to this problem are scattered across several areas of study. One of the contributions of this work is an organized summary of efforts in rational l_2 optimization, particularly related to the design of IIR

digital filters. The work in Section III-B also lays down a framework for the IIR methods proposed in this work.

As mentioned in Section III-C, some complications arise when designing IIR l_p filters. Aside from the intrinsic l_2 problem, it is necessary to properly combine a number of ideas that allowed for robust and efficient l_p FIR methods. A design algorithm for complex l_p IIR filters were presented in Section III-C1; this algorithm combined Soewito's quasilinearization with ideas such as l_p homotopy, partial updating and the adaptive modification. In practice, the combination of these ideas showed to be practical and the resulting algorithm remained robust. It was also found that after a few p -steps, the internal l_2 algorithm required from one to merely a few iterations on average, thus maintaining the algorithm efficient.

One of the main contributions of this work is the introduction of an IRLS-based method to solve l_p IIR design problems. By properly combining the principle of magnitude approximation via phase updating (from Soewito, Jackson and Kay) with the complex IIR algorithm one can find optimal magnitude l_p designs. This work also introduced a sequence of steps that improve the efficiency and robustness of this algorithm, by dividing the design process into three stages and by using suitable initial guesses for each stage.

Some of the examples in this document were designed using Matlab programs. It is worth to notice the common elements between these programs, alluding to the modularity of the implementations. An added benefit to this setup is that further advances in any of the topics covered in this work can easily be ported to most if not all of the algorithms.

Digital filter design is and will remain an important topic in digital signal processing. It is the hope of the author to have motivated in the reader some curiosity for the use of l_p norms as design criteria for applications in FIR and IIR filter design. This work is by no means comprehensive, and is meant to inspire the consideration of the flexibility of IRLS algorithms for new l_p related problems.

REFERENCES

- [1] R. E. Ziemer, W. H. Tranter, and D. R. Fannin, *Signals and Systems: Continuous and Discrete*, 4th ed. Prentice Hall, 1998.
- [2] J. L. Walsh and T. S. Motzkin, "Polynomials of Best Approximation on an Interval," *Proceedings of the National Academy of Sciences, USA*, vol. 45, pp. 1523–1528, October 1959.
- [3] T. S. Motzkin and J. L. Walsh, "Polynomials of Best Approximation on a Real Finite Point Set I," *Trans. American Mathematical Society*, vol. 91, no. 2, pp. 231–245, May 1959.
- [4] C. L. Lawson, "Contributions to the theory of linear least maximum approximations," Ph.D. dissertation, UCLA, 1961.
- [5] J. R. Rice and K. H. Usow, "The Lawson Algorithm and Extensions," *Mathematics of Computation*, vol. 22, pp. 118–127, 1968.
- [6] J. R. Rice, *The Approximation of Functions*. Addison-Wesley, 1964, vol. 1.
- [7] C. S. Burrus, J. A. Barreto, and I. W. Selesnick, "Iterative Reweighted Least-Squares Design of FIR Filters," *IEEE Transactions on Signal Processing*, vol. 42, no. 11, pp. 2926–2936, November 1994.
- [8] L. A. Karlovitz, "Construction of Nearest Points in the L^p , p even and L^∞ norms, I," *Journal of Approximation Theory*, vol. 3, pp. 123–127, 1970.
- [9] S. W. Kahng, "Best L_p Approximations," *Mathematics of Computation*, vol. 26, no. 118, pp. 505–508, April 1972.
- [10] R. Fletcher, J. A. Grant, and M. D. Hebden, "The Calculation of Linear Best L_p Approximations," *The Computer Journal*, vol. 14, no. 118, pp. 276–279, Apr 1972.
- [11] M. Aoki, *Introduction to Optimization Techniques*. The Macmillan Company, 1971.
- [12] J. A. Barreto, " L_p Approximation by the Iterative Reweighted Least Squares Method and the Design of Digital FIR Filters in One Dimension," Master's thesis, Rice University, 1992.
- [13] C. S. Burrus and J. A. Barreto, "Least p -power Error Design of FIR Filters," in *Proc. IEEE Int. Symp. Circuits, Syst. ISCAS-92*, vol. 2, San Diego, CA, May 1992, pp. 545–548.
- [14] J. Nocedal and S. J. Wright, *Numerical Optimization*, ser. Springer series in operations research. New York, NY: Springer-Verlag, 1999.
- [15] J. W. Adams and J. L. Sullivan, "Peak-Constrained Least-Squares Optimization," *IEEE Trans. on Signal Processing*, vol. 46, no. 2, pp. 306–321, Febr. 1998.
- [16] J. W. Adams, "FIR Digital Filters with Least-Squares Stopbands Subject to Peak-Gain Constraints," *IEEE Trans. on Circuits and Systems*, vol. 39, no. 4, pp. 376–388, April 1991.
- [17] I. W. Selesnick, M. Lang, and C. S. Burrus, "Constrained Least Square Design of FIR Filters without Specified Transition Bands," *IEEE Transactions on Signal Processing*, vol. 44, no. 8, pp. 1879–1892, August 1996.
- [18] M. Lang, I. W. Selesnick, and C. S. Burrus, "Constrained Least Square Design of 2-D FIR Filters," *IEEE Transactions on Signal Processing*, vol. 44, no. 5, pp. 1234–1241, May 1996.
- [19] I. W. Selesnick, M. Lang, and C. S. Burrus, "A Modified Algorithm for Constrained Least Square Design of Multiband FIR Filters Without Specified Transition Bands," *IEEE Transactions on Signal Processing*, vol. 46, no. 2, pp. 497–501, Feb. 1998.
- [20] R. Fletcher and M. J. D. Powell, "A Rapidly Convergent Descent Method for Minimization," *Computer Journal*, vol. 6, no. 2, pp. 163–168, 1963.
- [21] J. E. Dennis and R. B. Schnabel, *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*. Philadelphia, PA: SIAM, 1996.
- [22] L. S. e. a. T. P. Krauss, *Signal Processing Toolbox User's Guide*. The MathWorks, 1994, ch. 2, pp. 143–145.
- [23] J. T. Spanos and D. L. Mingori, "Newton Algorithm for Fitting Transfer Functions to Frequency Measurements," *Journal of Guidance, Control and Dynamics*, vol. 16, no. 1, pp. 34–39, January 1993.
- [24] D. C. Sorensen, "Newton's Method with a Dodel Trust Region Modification," *SIAM Journal of Numerical Analysis*, vol. 16, pp. 409–426, 1982.
- [25] R. Kumaresan and C. S. Burrus, "Fitting a Pole-Zero Filter Model to Arbitrary Frequency Response Samples," *Proc. ASIOMAR*, pp. 1649–1652, 1991.
- [26] F. B. Hildebrand, *Introduction to Numerical Analysis*. McGraw-Hill, 1974.
- [27] E. C. Levy, "Complex-Curve Fitting," *IRE Transactions on Automatic Control*, vol. AC-4, no. 1, pp. 37–43, May 1959.
- [28] T. W. Parks and C. S. Burrus, *Digital Filter Design*. John Wiley and Sons, 1987.
- [29] B. G. C. F. M. R. de Prony, "Essai Expérimental et Analytique: Sur les lois de la Dilatabilité des fluides élastiques et sur celles de la Force expansive de la vapeur de l'eau et de la vapeur de l'alcool, à différentes températures," *Journal de l'École Polytechnique (Paris)*, vol. 1, no. 2, pp. 24–76, 1795.
- [30] H. E. Padé, "Sur la Représentation Approchée d'une Fonction par des Fractions Rationnelles," *Annales Scientifiques de L'École Normale Supérieure (Paris)*, vol. 9, no. 3, pp. 1–98, 1892.
- [31] C. K. Sanathanan and J. Koerner, "Transfer Function Synthesis as a Ratio of Two Complex Polynomials," *IEEE Transactions on Automatic Control*, vol. AC-8, pp. 56–58, January 1963.
- [32] A. C. M. A. Sid-Ahmed and G. A. Jullien, "Computational Techniques for Least-Square Design of Recursive Digital Filters," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. ASSP-26, no. 5, pp. 477–480, October 1978.
- [33] H. W. S. L. E. McBride and K. Steiglitz, "Time-Domain Approximation by Iterative Methods," *IEEE Transactions on Circuit Theory*, vol. CT-13, no. 4, pp. 381–87, December 1966.
- [34] K. Steiglitz and L. E. McBride, "A Technique for the Identification of Linear Systems," *IEEE Transactions on Automatic Control*, vol. AC-10, pp. 461–64, October 1965.
- [35] A. W. Soewito, "Least square digital filter design in the frequency domain," Ph.D. dissertation, Rice University, December 1990.
- [36] L. B. Jackson, "Frequency-Domain Steiglitz-McBride Method for Least-Squares IIR Filter Design, ARMA Modeling, and Periodogram Smoothing," *IEEE Signal Processing Letters*, vol. 15, pp. 49–52, 2008.
- [37] A. Antoniou, *Digital Filters: Analysis, Design, and Applications*, 2nd ed. McGraw-Hill, 1993.
- [38] E. Cunningham, *Digital Filtering: An Introduction*. Houghton-Mifflin, 1992.
- [39] E. W. Cheney, *Introduction to Approximation Theory*, ser. Intl. Series in Pure and Applied Mathematics. McGraw-Hill, 1966.
- [40] V. Chvatal, *Linear Programming*. Freeman and Co., 1980.
- [41] G. Strang, *Introduction to Applied Mathematics*. Wellesley-Cambridge Press, 1986.

- [42] S. A. Ruzinsky, " L_1 and L_∞ Minimization via a Variant of Karmarkar's Algorithm," *IEEE Trans. on Acoustics, Speech and Signal Processing*, vol. 37, no. 2, pp. 245–253, Febr. 1989.
- [43] A. V. Oppenheim and R. W. Schaffer, *Discrete-Time Signal Processing*. Englewood Cliffs, N.J.: Prentice-Hall, 1989.
- [44] C. S. Burrus, J. H. McClellan *et al.*, *Computer-based Exercises for Signal Processing*. Prentice Hall, 1994.
- [45] R. A. Vargas and C. S. Burrus, "Adaptive Iterative Reweighted Least Squares Design of L_p FIR Filters," *Proc. ICASSP*, vol. 3, Signal Processing, Theory and Methods, pp. 1129–32, 1999.
- [46] K. S. Shanmugan and A. M. Breipohl, *Random Signals: Detection, Estimation and Data Analysis*. John Wiley & Sons, 1988.
- [47] J. G. Proakis and D. G. Manolakis, *Digital Signal Processing*. Macmillan Publishing Co., 1988.
- [48] S.-P. Wu, S. Boyd, and L. Vandenberghe, "Fir filter design via spectral factorization and convex optimization," to Appear in *Applied Computational Control, Signal and Communications*, Biswa Datt editor, Birkhauser, 1997.
- [49] K. Steiglitz, "Computer-Aided Design of Recursive Digital Filters," *IEEE Transactions on Audio and Electroacoustics*, vol. AU-18, no. 2, pp. 123–129, June 1970.
- [50] A. Deczky, "Synthesis of Recursive Digital Filters Using the Minimum p -error Criterion," *IEEE Transactions on Audio and Electroacoustics*, vol. AU-20, no. 4, pp. 257–263, October 1972.
- [51] R. Kumaresan, "Identification of Rational Transfer Functions from Frequency Response Samples," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 26, no. 6, pp. 925–934, November 1990.
- [52] R. E. Mickens, *Difference equations*. Van Nostrand Reinhold, 1987.
- [53] S. N. Elaydi, *An Introduction to Difference Equations*, ser. Undergraduate texts in mathematics. New York: Springer, 1996.
- [54] D. F. T. Jr., "On Fluids, Networks, and Engineering Education," in *Aspects of Network*, R. E. Kalman and N. DeClaris, Eds. Hol, Rinehart and Winston, Inc, 1971, pp. 591–612.
- [55] M. L. V. Blaricum, "A Review of Prony's Method techniques for Parameter Estimation," in *Air Force Statistical Estimation Workshop*, May 1978, pp. 125–135.
- [56] L. Weiss and R. N. McDonough, "Prony's method, Z-transforms, and Padé approximation," *SIAM Review*, vol. 5, no. 2, pp. 145–149, April 1963.
- [57] I. Barrondale and D. D. Olesky, "Exponential Approximation using Prony's Method," in *The Numerical Solution of Nonlinear Problems*, C. T. H. Baker and C. Phillips, Eds. New York: Oxford University Press, 1981, ch. 20, pp. 258–269.
- [58] S. L. Marple, Jr., *Digital Spectral Analysis with Applications*, ser. Signal Processing Series. Englewood Cliffs, NJ: Prentice-Hall, 1987.
- [59] C. F. Gerald and P. O. Wheatley, *Applied Numerical Analysis*. Reading, MA: Addison-Wesley, 1984.
- [60] H. Cabannes, Ed., *Pade Approximates Method and its Applications to Mechanics*, ser. Lecture notes in physics. Springer-Verlag, 1976, no. 47.
- [61] G. H. Golub and C. F. V. Loan, *Matrix Computations*. Johns Hopkins University Press, 1996.
- [62] H. Fan and M. Doroslovački, "On "Global Convergence" of Steiglitz-McBride Adaptive Algorithm," *IEEE Transactions on Circuits and Systems II*, vol. 40, no. 2, pp. 73–87, February 1993.
- [63] P. Stoica and T. Söderström, "The Steiglitz-McBride Identification Algorithm Revisited-Convergence Analysis and Accuracy Aspects," *IEEE Transactions on Automatic Control*, vol. AC-26, no. 3, pp. 712–17, June 1981.
- [64] J. H. McClellan and D. Lee, "Exact Equivalence of the Steiglitz-McBride Iteration and IQML," *IEEE Transactions on Signal Processing*, vol. 39, no. 2, pp. 509–12, February 1991.
- [65] R. E. Bellman and R. E. Kalaba, *Quasilinearization and Nonlinear Boundary-Value Problems*. New York: American Elsevier, 1965.
- [66] R. A. Vargas and C. S. Burrus, "On the Design of L_p IIR Filters with Arbitrary Frequency Bands," in *Proc. ICASSP*, vol. 6, 2001, pp. 3829–3832.